

Développement d'outils biostatistiques et bioinformatiques de prédiction et d'analyse des défauts de l'épissage : application aux gènes de prédisposition aux cancers du sein et de l'ovaire

Raphaël Leman

Laboratoire de biologie et génétique du cancer

Inserm U1245 Genomics and Personalized Medecine in Cancer and Neurological Disorders

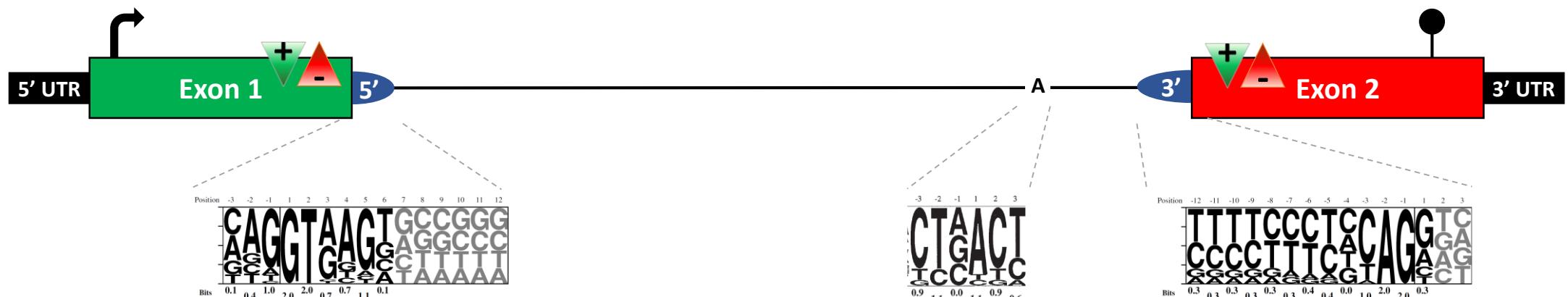
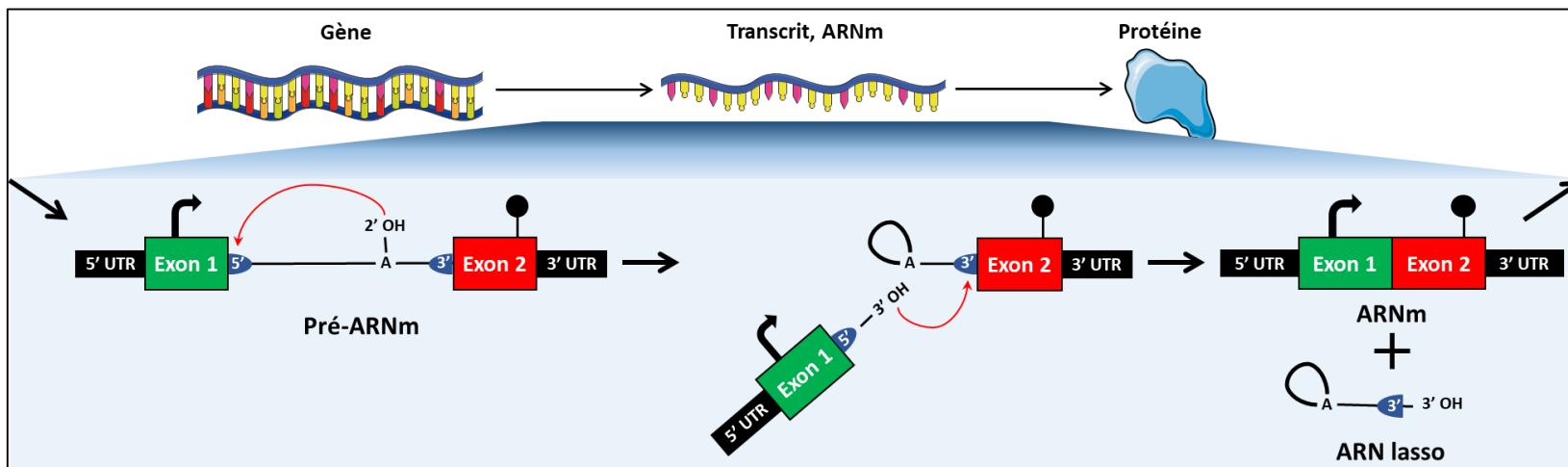
Centre François Baclesse

Normandie Univ, UNICAEN, Caen

**Encadré par : Dr Sophie Krieger
Co-Encadré par : Dr Alexandra Martins**

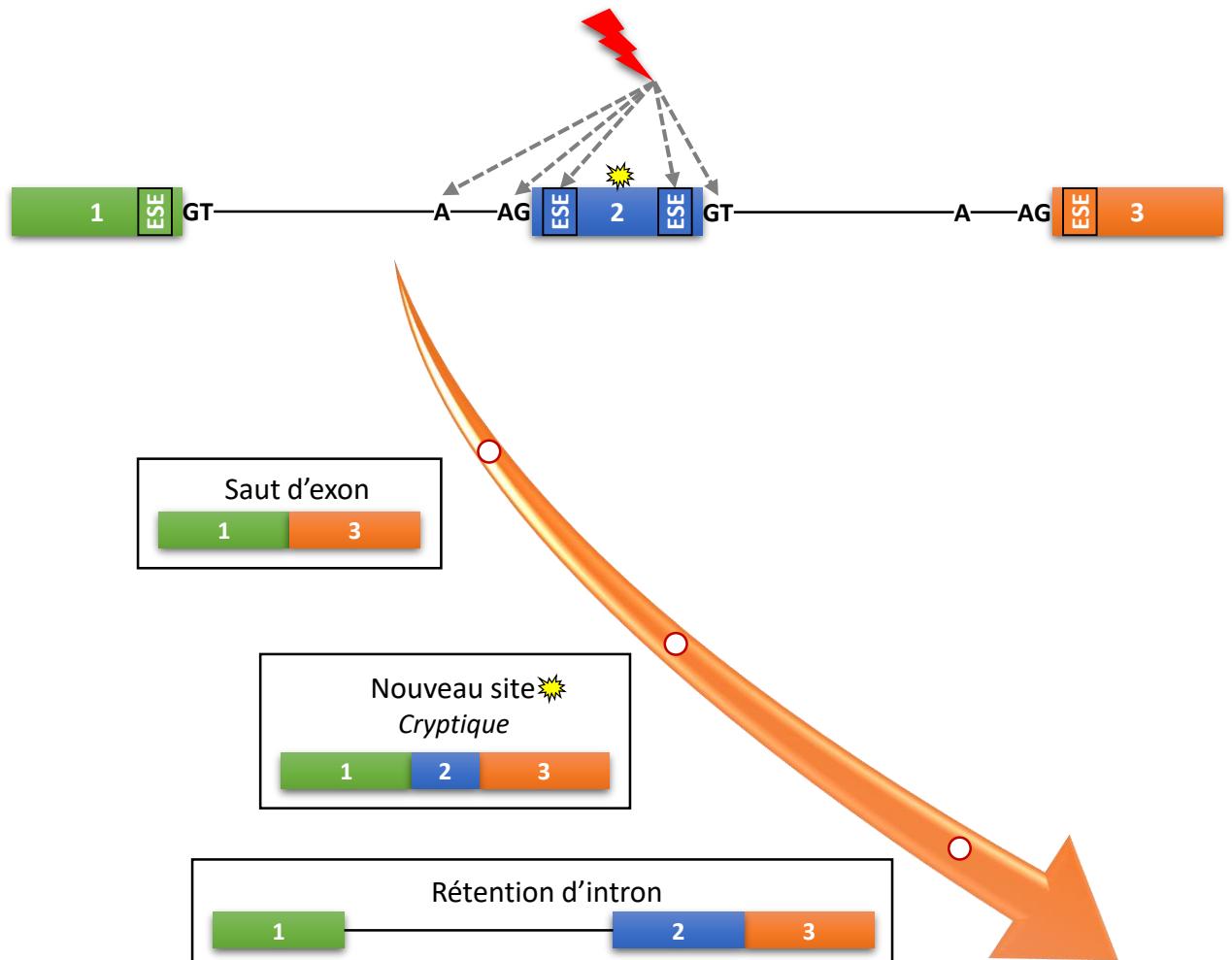


Les motifs d'épissage

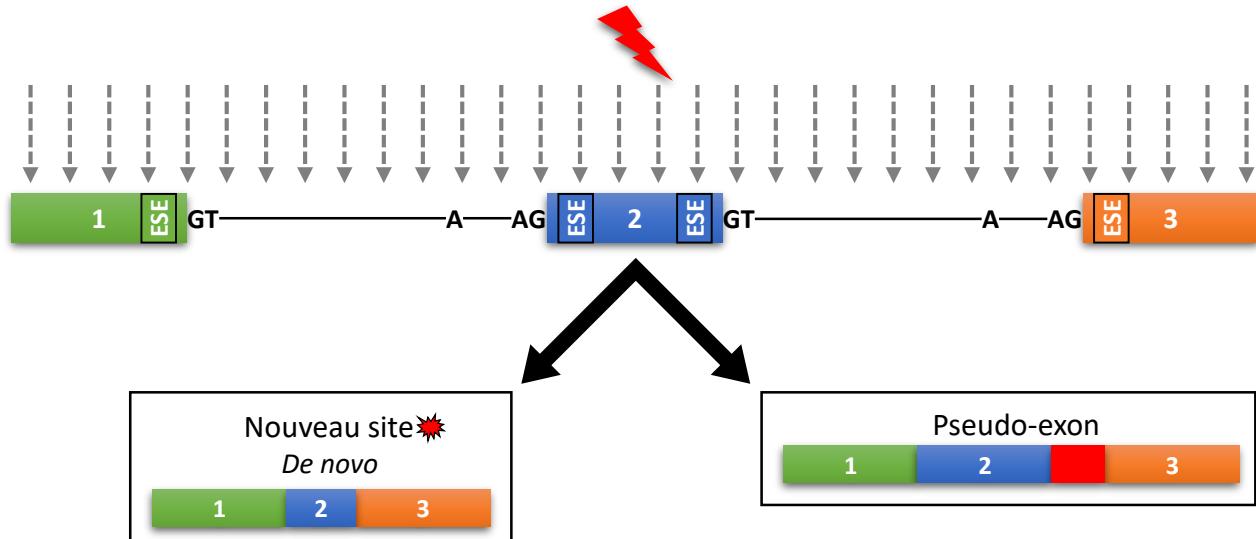


Mécanismes d'altération de l'épissage

Altération des motifs d'épissage

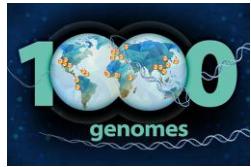


Création de motifs d'épissage



3,8 % des variations génétiques exoniques et introniques proches des exons **impactent l'épissage** en dehors de tout contexte de prédisposition

Comment identifier un variant splicéogénique ?



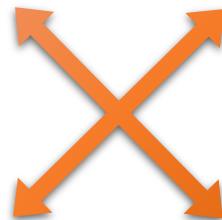
1000 Genomes Project
→ 2 500 genomes



84 000 000 variations génétiques



Population mondiale
→ 7 milliards



$2,4 \times 10^{14}$ variations génétiques ⇔ ~1 000 x

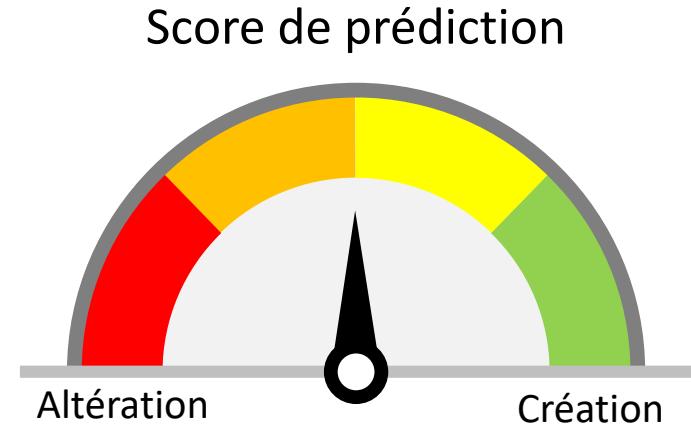


Voie lactée
→ 250 milliards d'étoiles

Impossible de tester expérimentalement chaque variant sur l'épissage

Prédire un variant splicéogénique

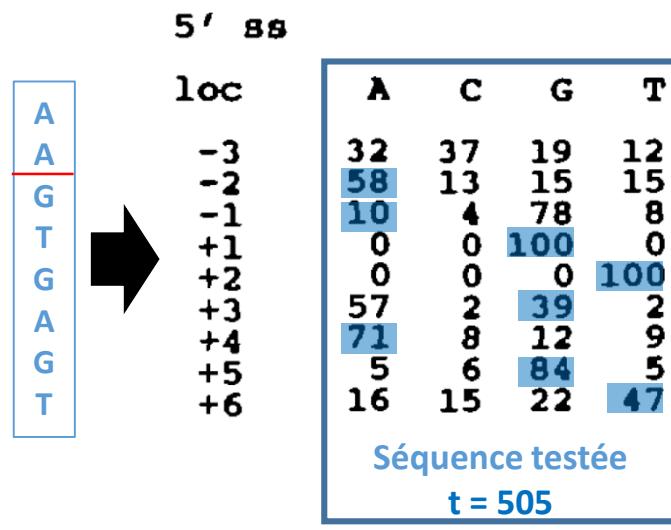
La prédition d'un variant splicéogénique \Leftrightarrow la prédition de la création/altération de motifs d'épissage



Un score par motif

Ex: Splice Site Finder (SSF)¹ et MaxEntScan (MES)²

SSF



Somme des % minimaux

Somme des % maximaux

$$Score_{SSF} = 100 \times \left(\frac{t - mint}{maxt - mint} \right)$$

$$Score_{SSF} = 84,3 \%$$

MES

Non dépendante

AGCTGATCAGTCTGATTCC
...
P(X₁) non dépendant de P(X₂), ...

De proche en proche

AGCTGATCAGTCTGATTCC
↑↑↑ ...
P_{X₁}(X₂) >0, P_{X₂}(X₃)>0, ...

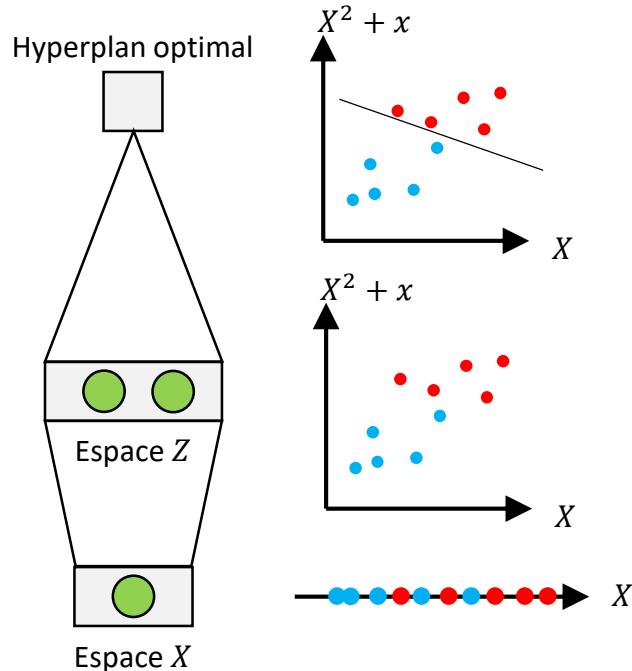
Discontinu

AGCTGATCAGTCTGATTCC
↑↑↑ ...
P_{X₂}(X₄) >0, P_{X₅}(X₉)>0, ...

Calcul à partir d'une séquence complète

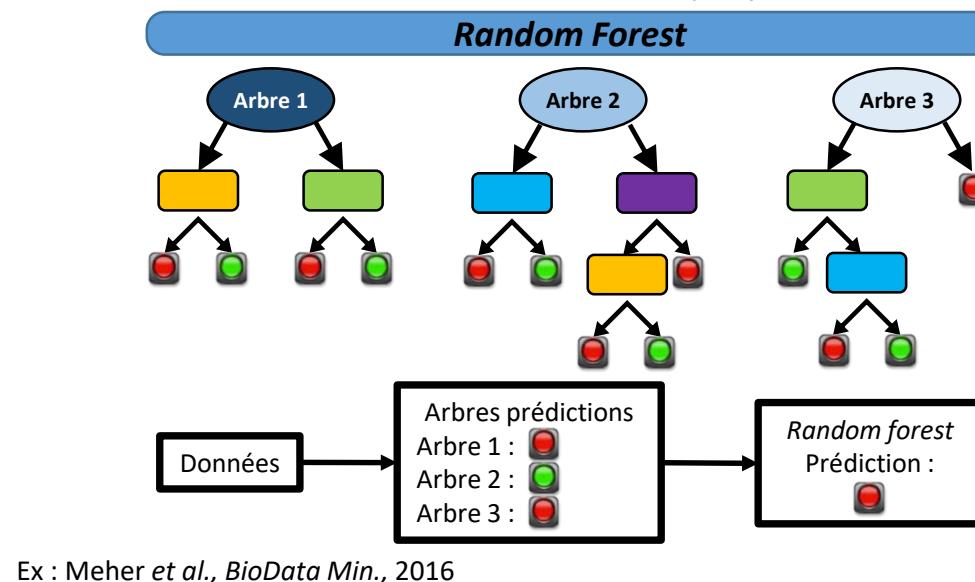
Grande diversité d'outils

Support Vector Machine (SVM)



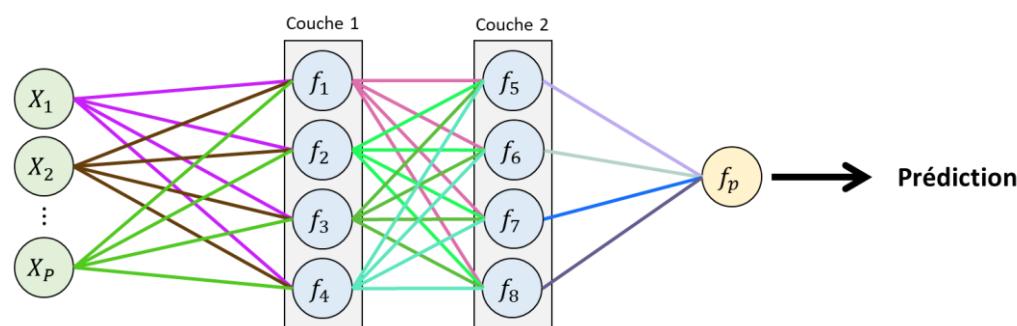
Ex : SVM-BPfinder; Corvelo, *PLOS Comput. Biol.*, 2010

Random Forest (RF)



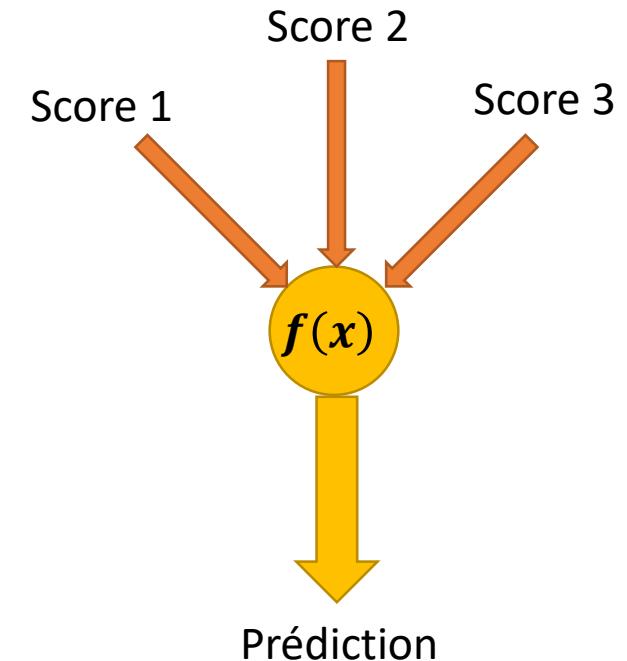
Ex : Meher et al., *BioData Min.*, 2016

Réseaux de neurones



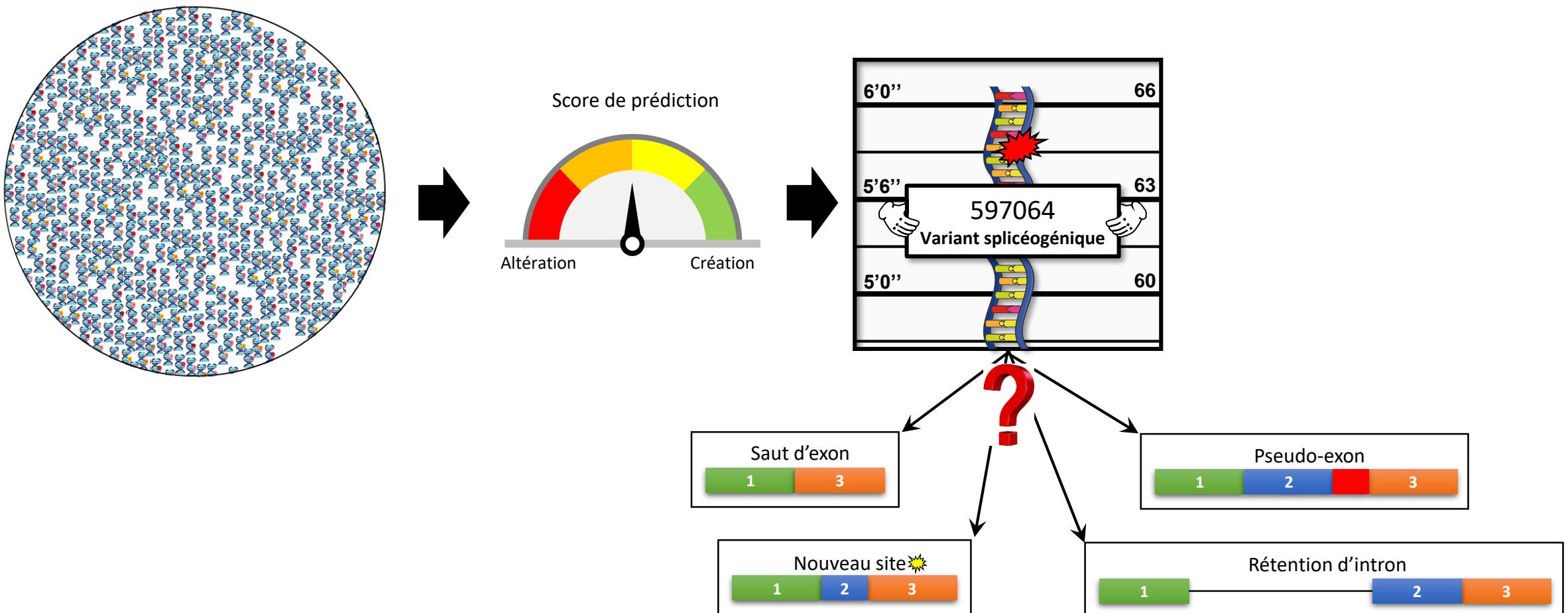
Ex : NNS (Neural Network Splice); Reese et al., *In Gene-Finding and Gene Structure*, 1995

Meta-Score



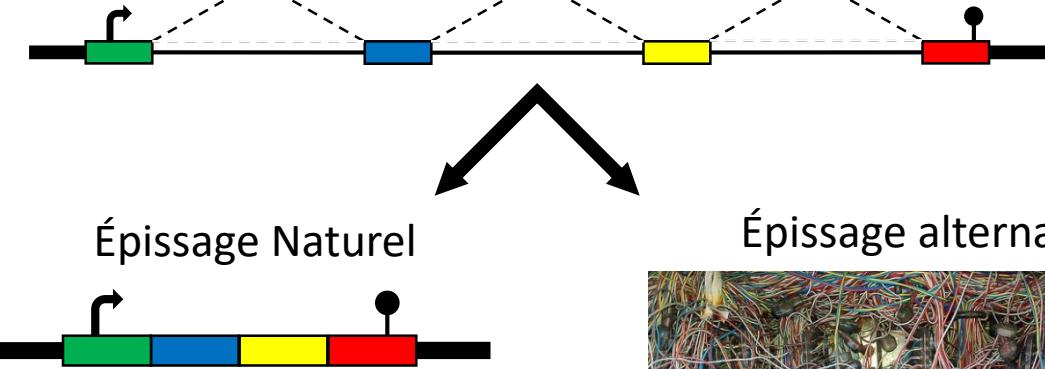
Ex : ADAboost; Jian et al., *Nucleic Acids Res.*, 2014

Étude des variants splicéogéniques



Épissage alternatif : rôle majeur dans la diversité des transcrits

À partir du même pré-mRNA



Épissage Naturel

Épissage alternatif



Saut d'exon

Utilisation nouveau site d'épissage

Rétention d'intron

10



Humain



Gène : *KCNMA1*
> 500 ARNm ≠

Drosophile



Gène : *Dscam*
> 35 000 ARNm ≠

Tests fonctionnels

Tests à haut débit

RNA-seq



1 Gb – 1 Tb

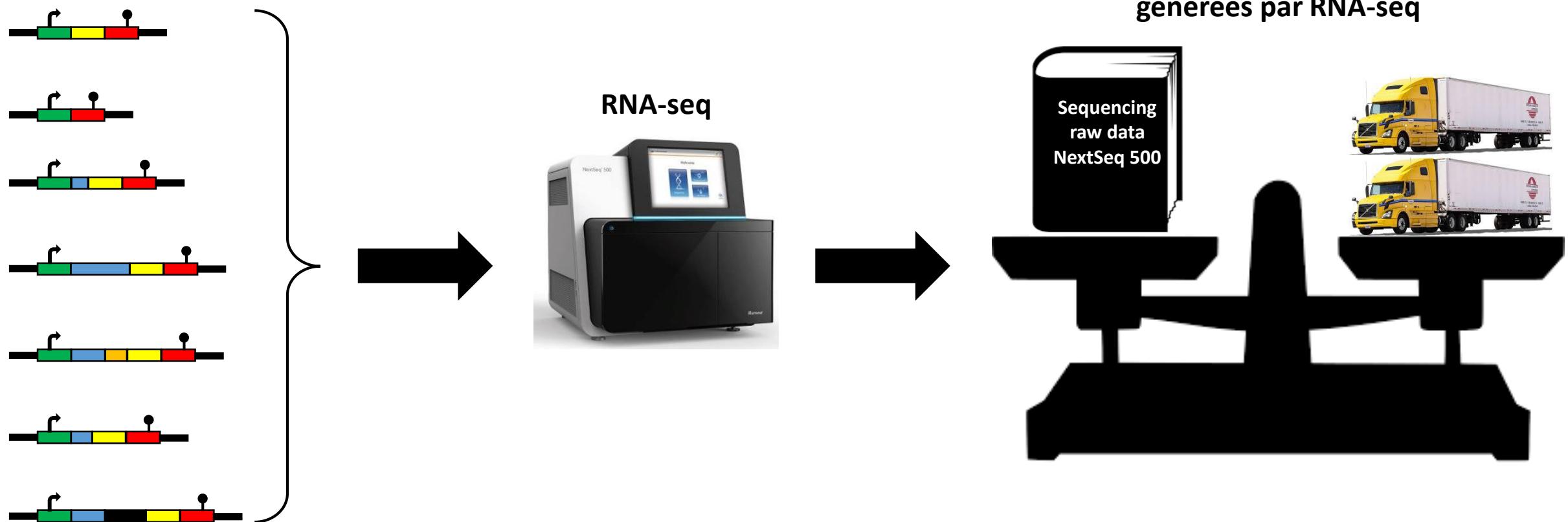
Tests à bas débit

RT-PCR, test
minigène,
electrophorese
capillaire, ...

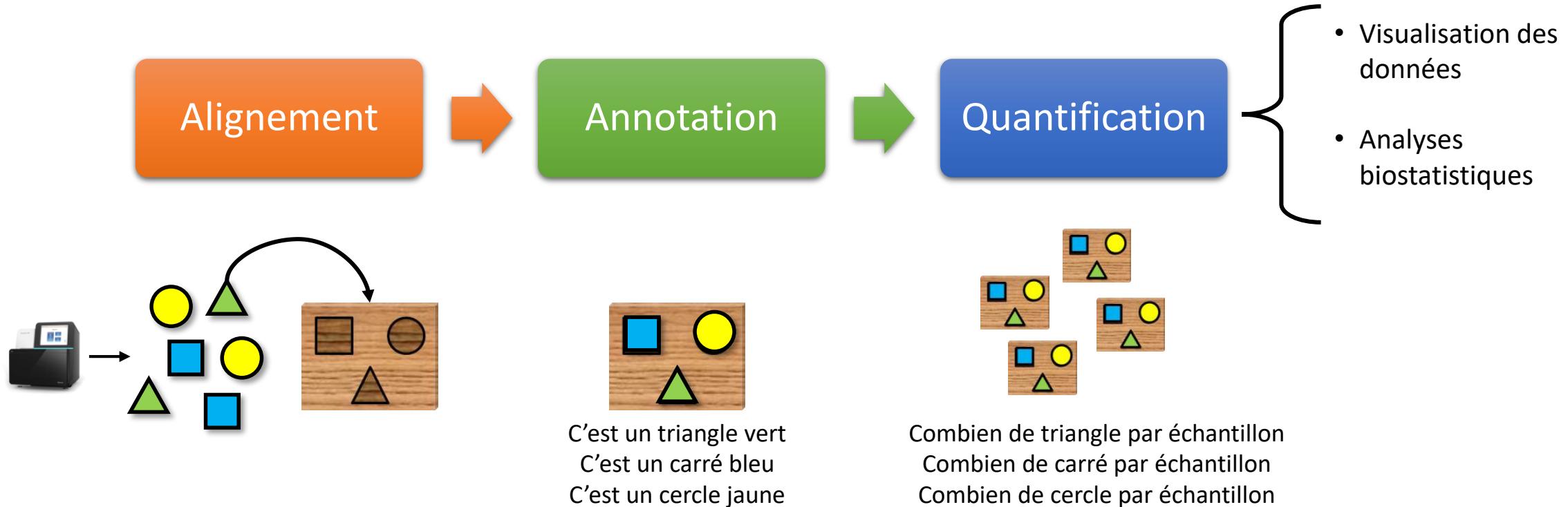


< 10 kb

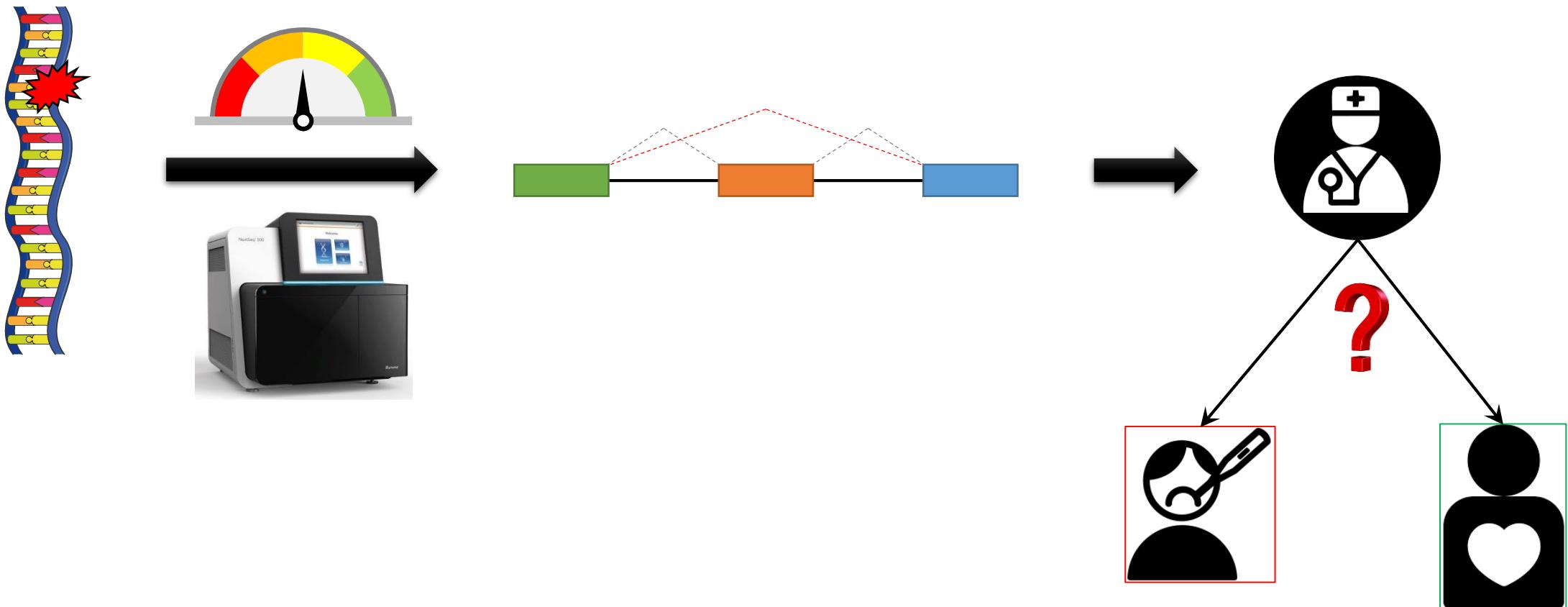
RNA-seq peut répondre à la diversité des transcrits
→ Comment extraire la structure de ces transcrits



RNA-seq et analyses bioinformatiques

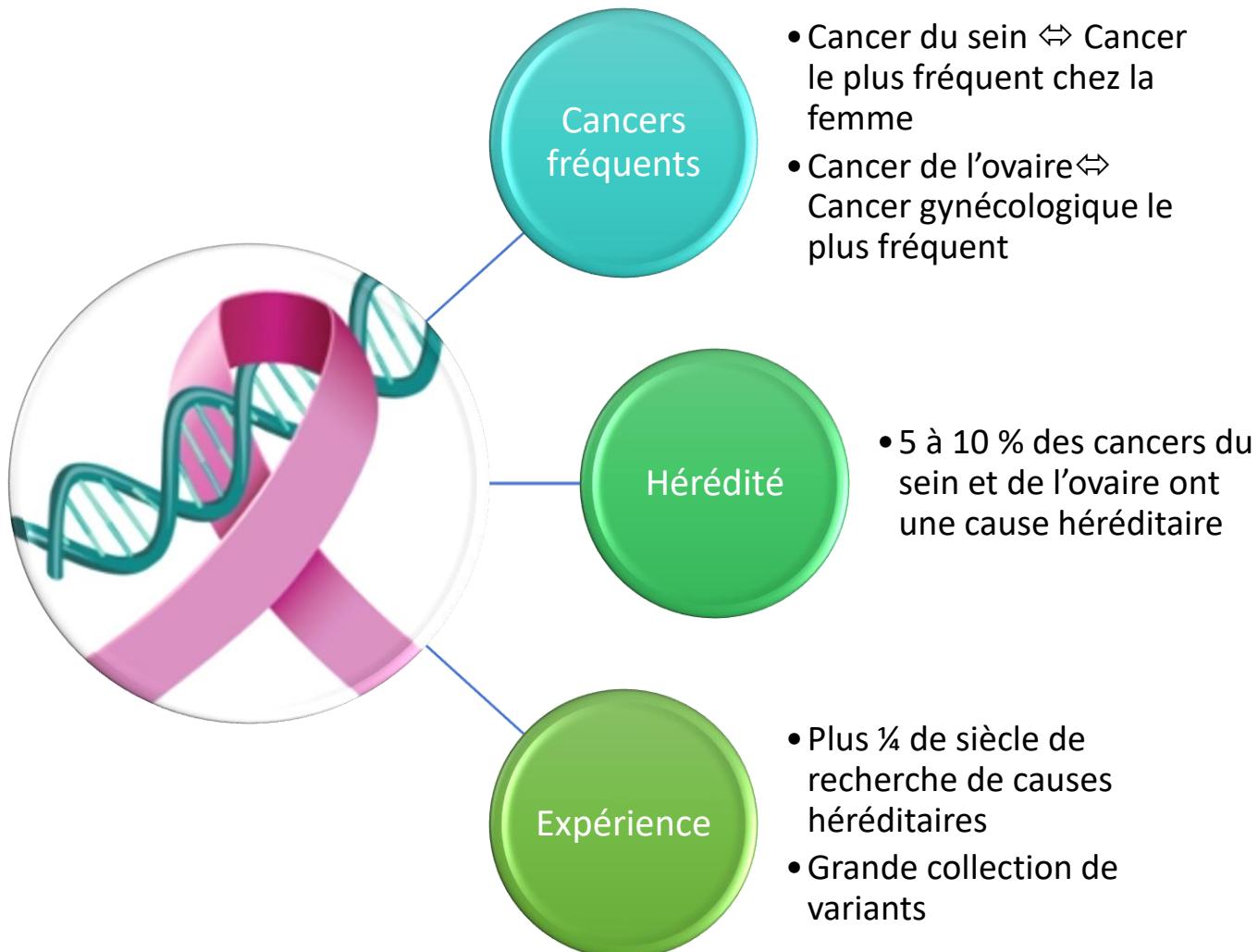


Transcrits et phénotypes



Épissage et pathogénicité

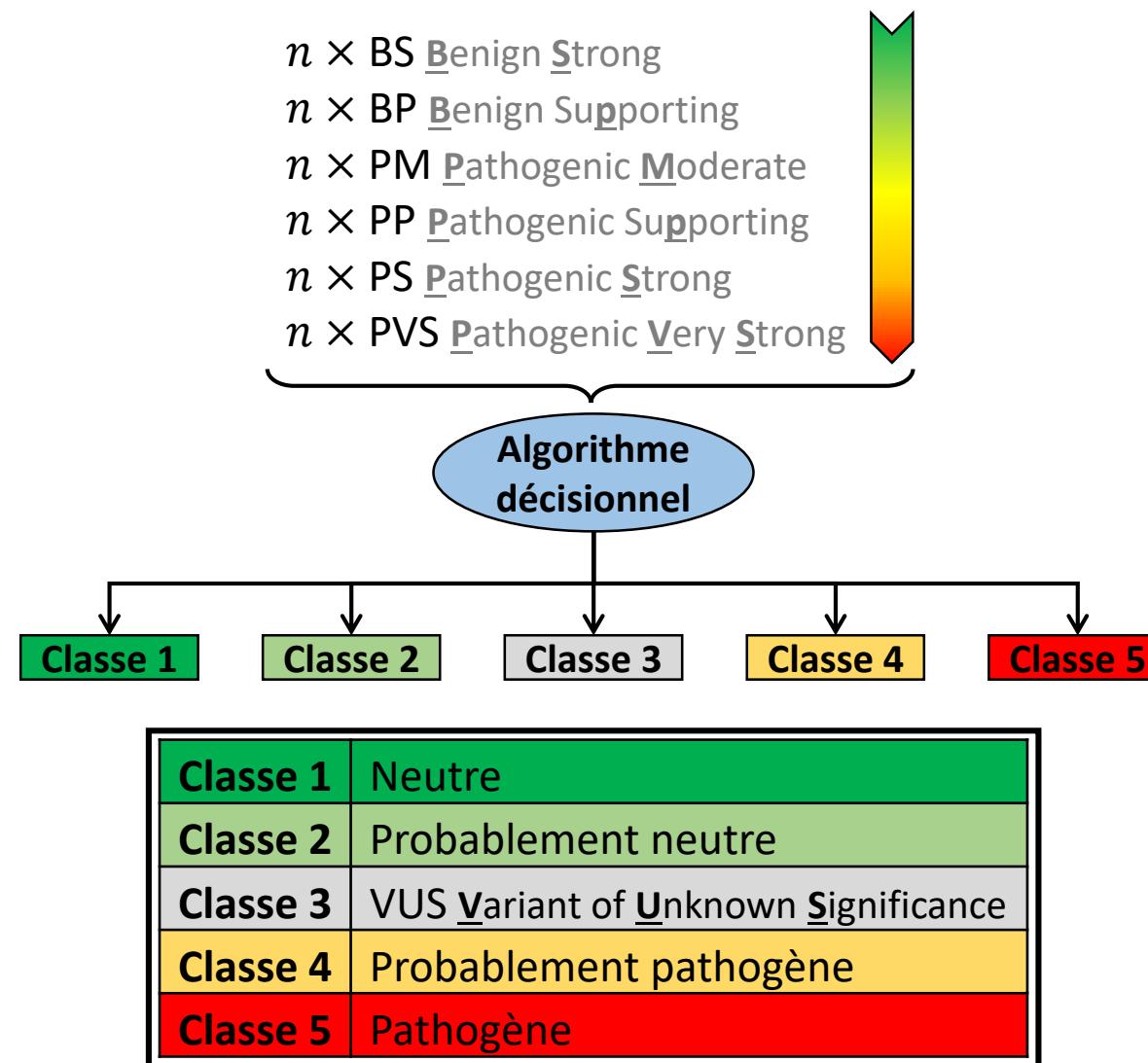
Étude de la prédisposition aux cancers du sein et de l'ovaire



5 gènes majeurs de prédisposition :

- *BRCA1* (1994)
- *BRCA2* (1995)
- *PALB2* (2007)
- *RAD51C* (2010)
- *RAD51D* (2011)

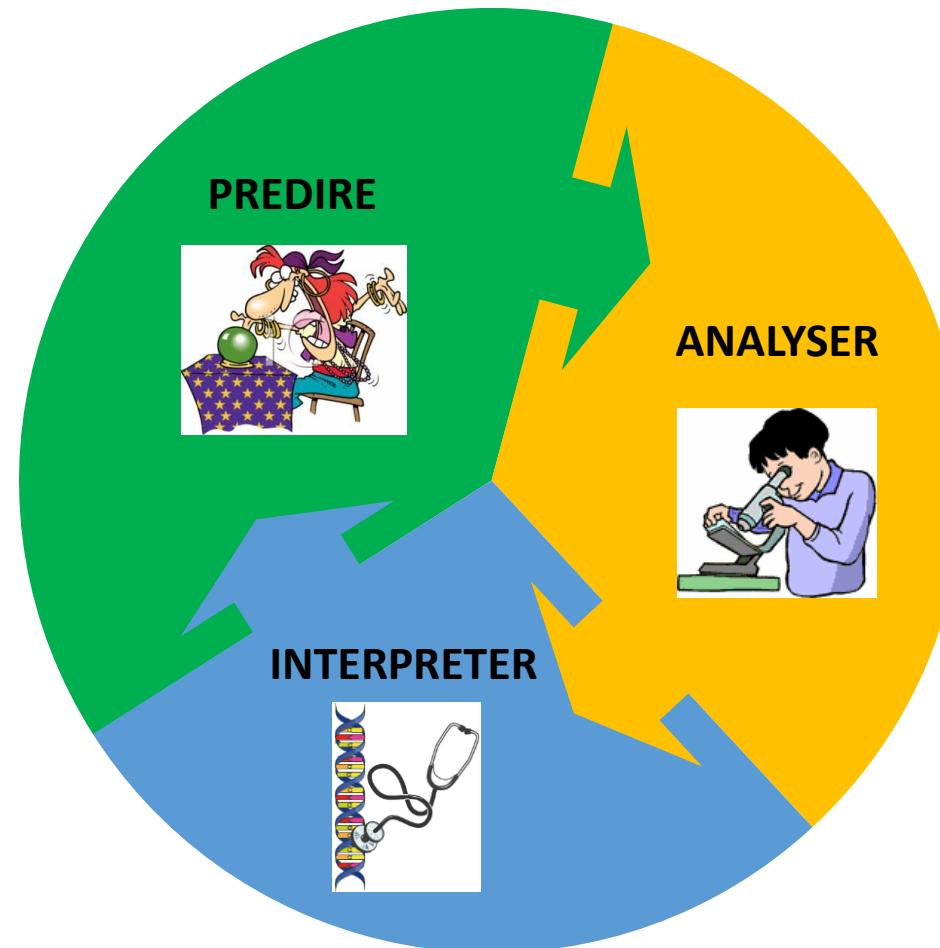
Classification des variants à usage clinique

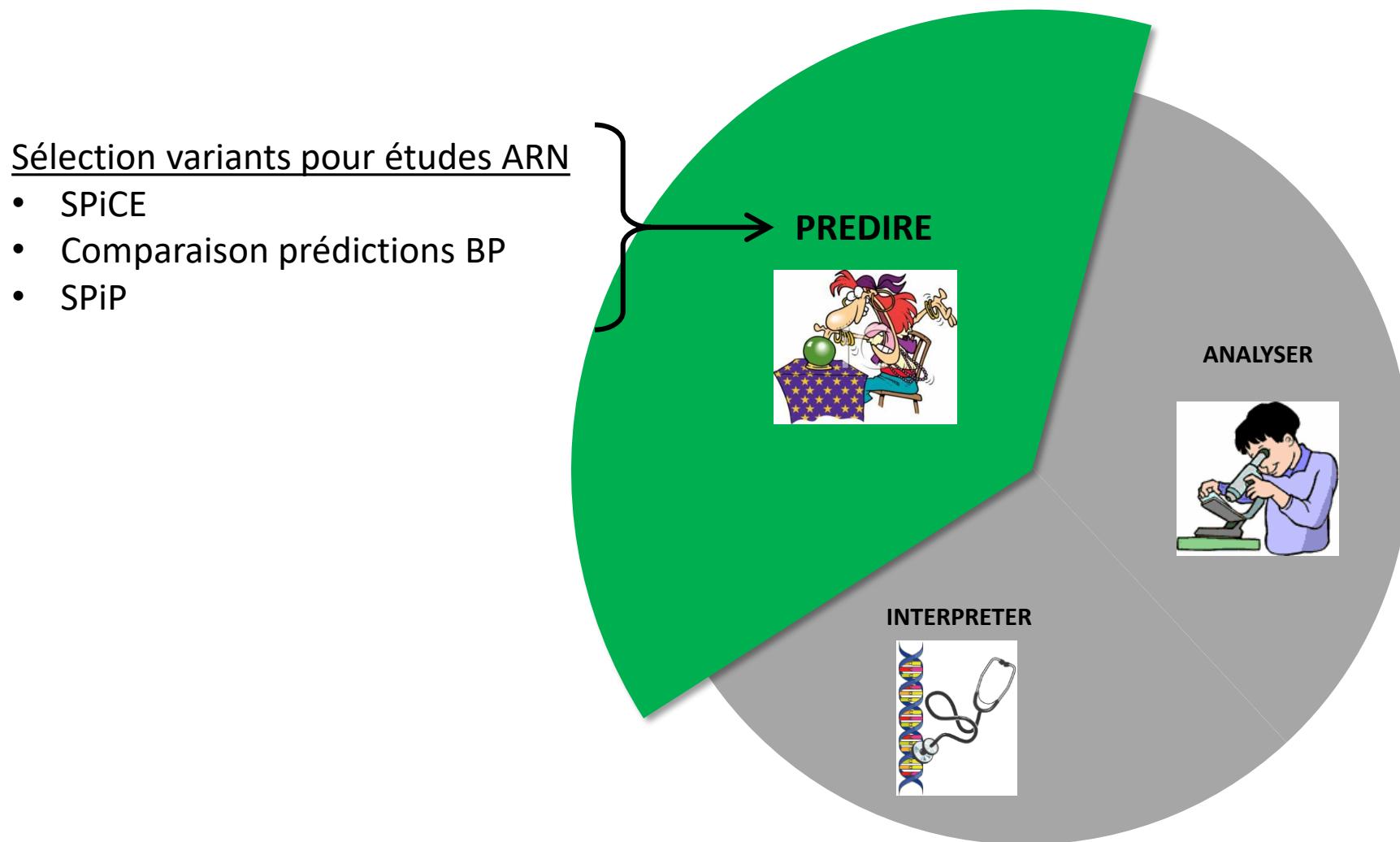


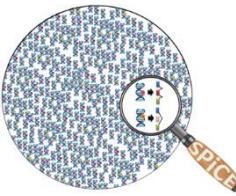
Les **PVS** comprennent entre autres les variants **tronquants** et **canoniques** (-1/-2; +1/+2)

Objectif des travaux de thèse

Le cycle des études des défauts d'épissage



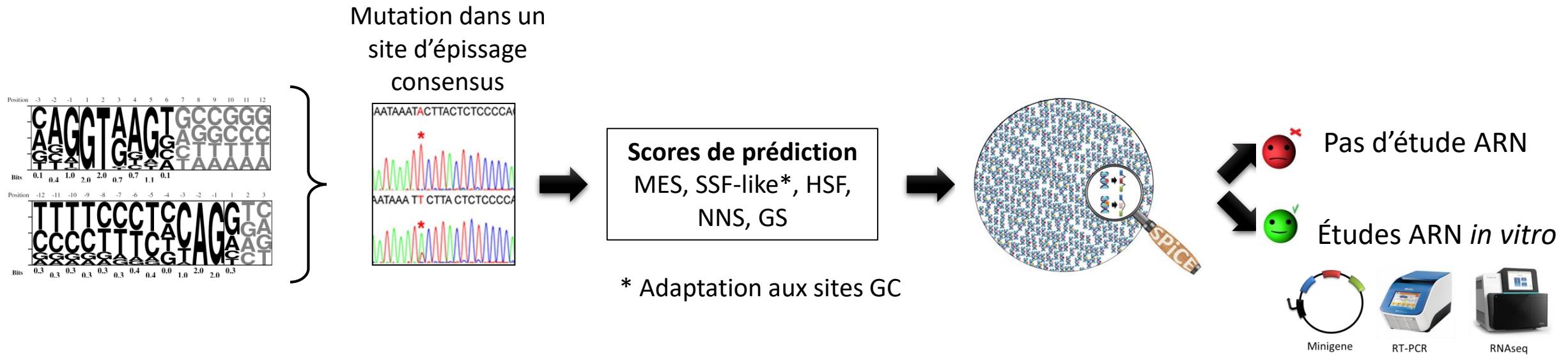


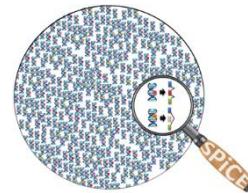


SPiCE, Splicing Prediction in Consensus Element (1/6)

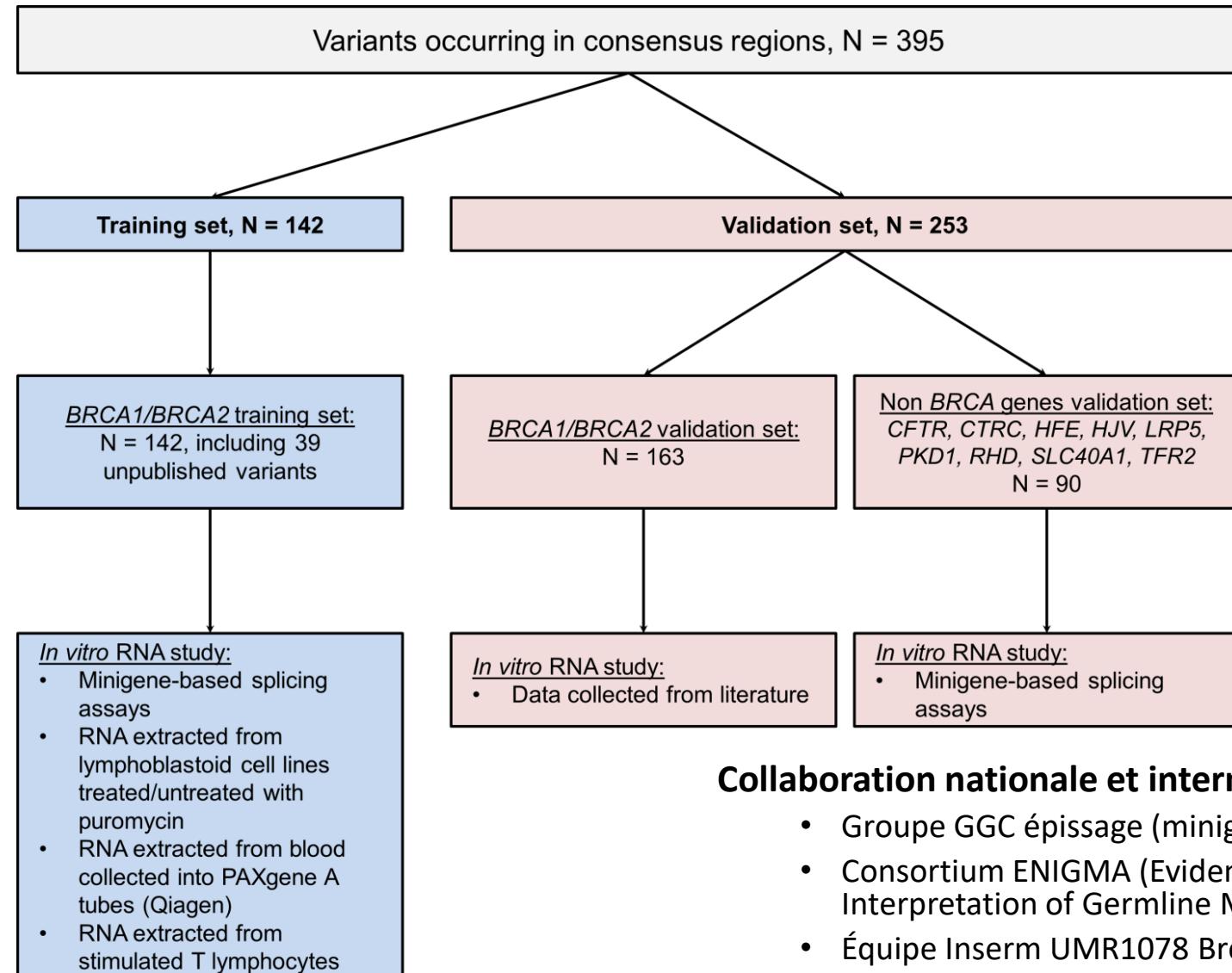
- **Objectif:**

À partir des prédictions bioinformatiques, être capable de prioriser l'étude ARN des variants à haut risque d'altérer l'épissage



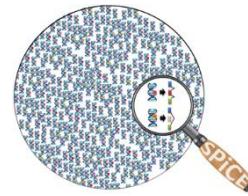


Jeux de données (2/6)



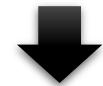
Collaboration nationale et internationale :

- Groupe GGC épissage (minigènes U1245, Rouen)
- Consortium ENIGMA (Evidence-based Network for the Interpretation of Germline Mutant Alleles)
- Équipe Inserm UMR1078 Brest



Construction de l'outil (3/6)

Données d'apprentissage
+
Régression logistique



Variables utilisées

- Scores de prédictions (SSF-like, MES, HSF, NNS et GS)



Sélection variables

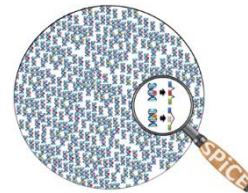


Validation modèle



Modèle final

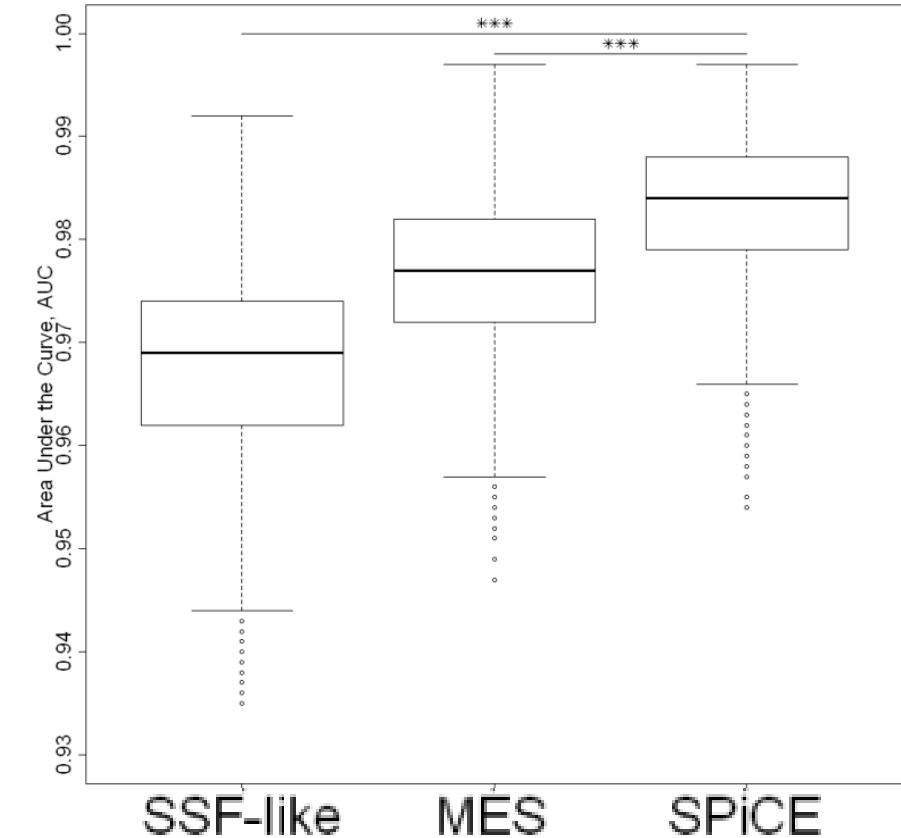
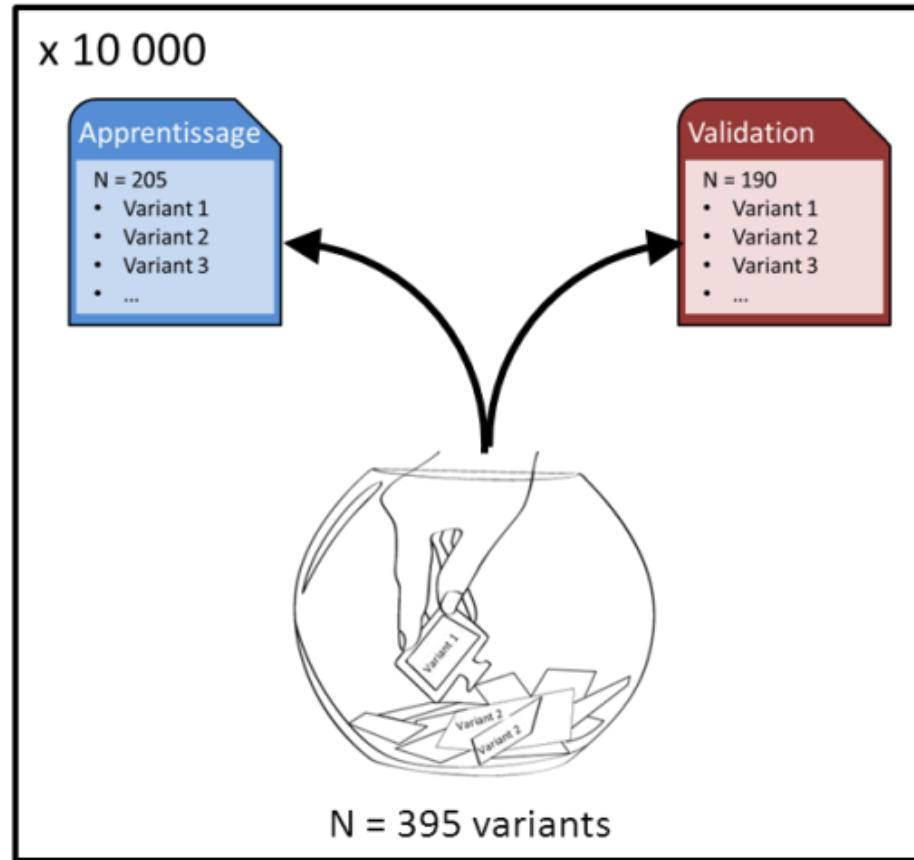
+ cross-validation

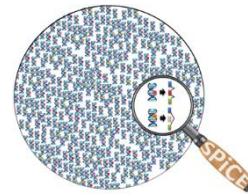


Validation de l'outil (4/6)

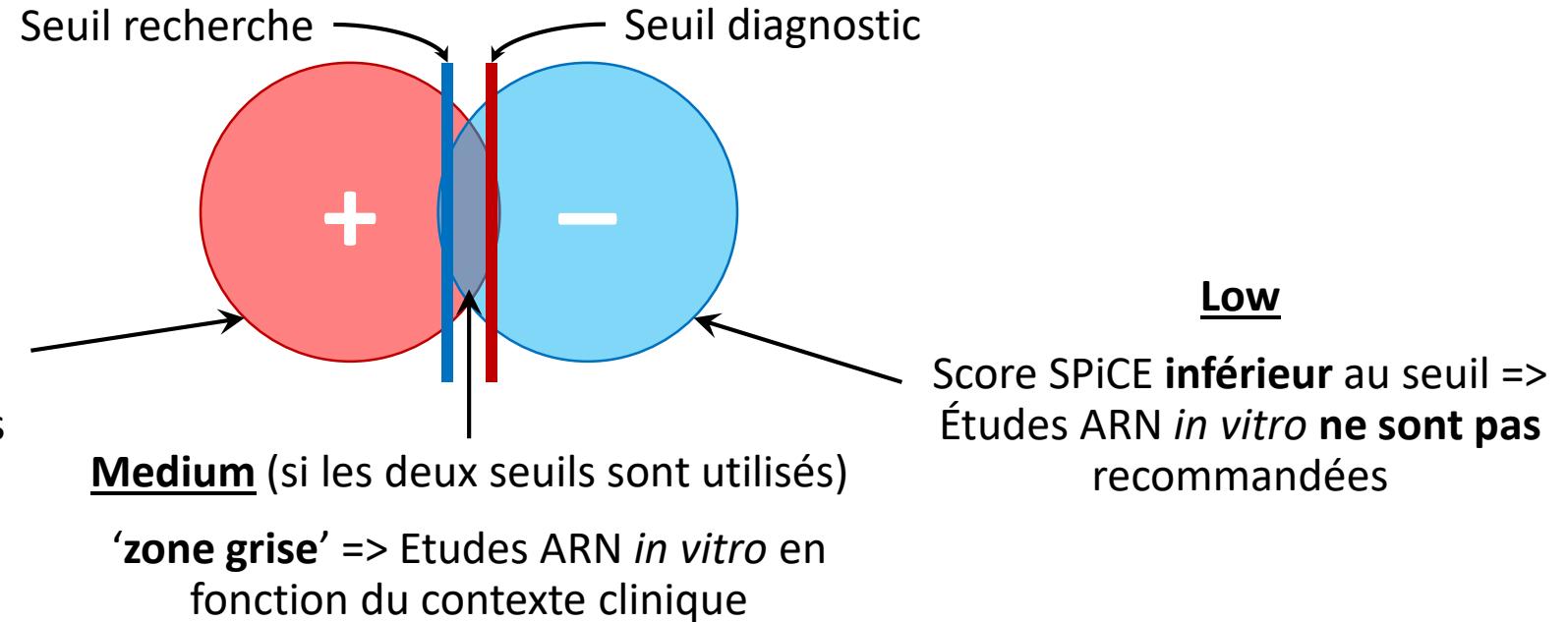
SPiCE : Combinaison par régression logistique
des deux scores **MES** et **SSF-like**

Cross-Validation





Performances et règles décisionnelles (5/6)



Seuil adapté au diagnostic moléculaire :

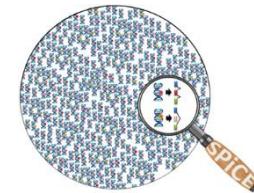
Sensibilité = 99,5 %

Spécificité = 75,6 %

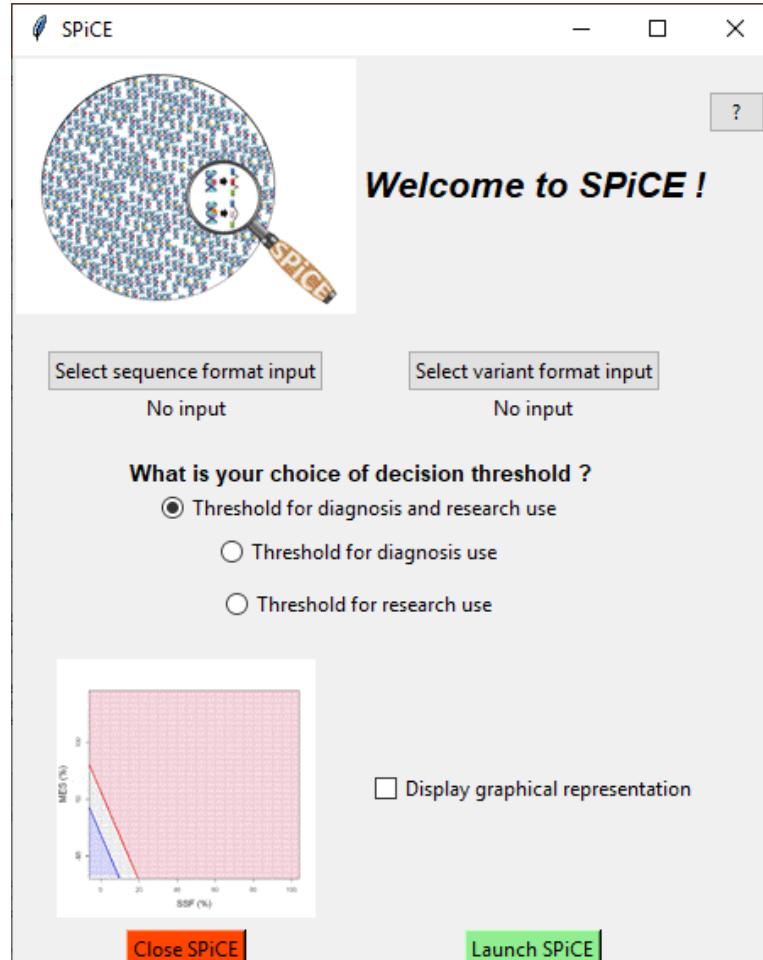
Seuil adapté à la recherche :

Sensibilité = 92,3 %

Spécificité = 95,6 %



Application et diffusion de SPiCE (6/6)



Novel diagnostic tool for prediction of variant spliceogenicity derived from a set of 395 combined *in silico/in vitro* studies: an international collaborative effort

Raphaël Leman, Pascaline Gaildrat, Gérald L Gac, Chandran Ka, Yann Fichou, Marie-Pierre Audrezet, Virginie Caux-Moncoutier, Sandrine M Caputo, Nadia Boutry-Kryza, Mélanie Léone, Sylvie Mazoyer, Françoise Bonnet-Dorion, Nicolas Sevenet, Marine Guillaud-Bataille, Etienne Rouleau, Brigitte Bressac-de Paillerets, Barbara Wappenschmidt, Maria Rossing, Danielle Muller, Violaine Bourdon, Françoise Revillon, Michael T Parsons, Antoine Rousselin, Grégoire Davy, Gaia Castelain, Laurent Castéra, Joanna Sokolowska, Florence Coulet, Capucine Delnatte, Claude Férec, Amanda B Spurdle, Alexandra Martins, Sophie Krieger , Claude Houdayer

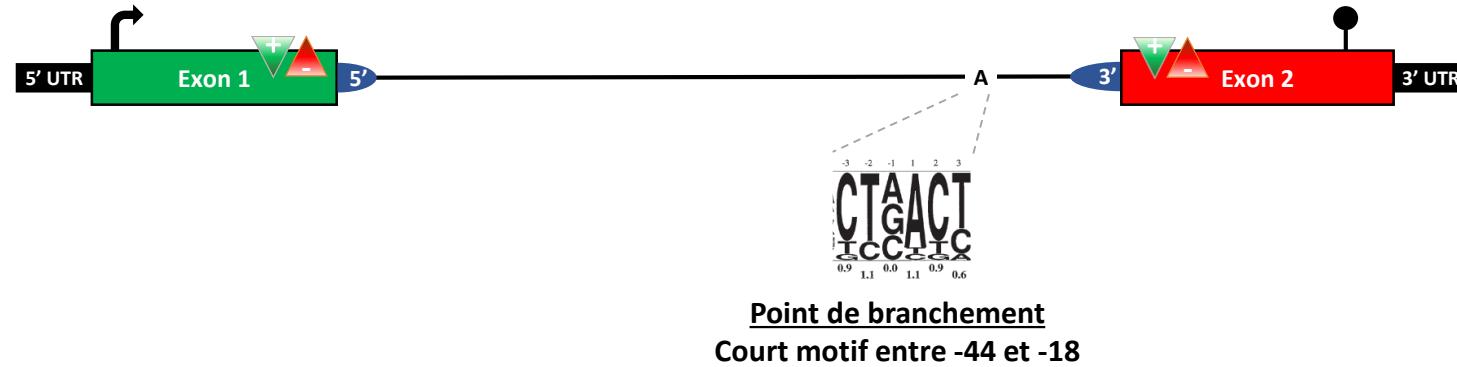
Author Notes

Nucleic Acids Research, Volume 46, Issue 15, 6 September 2018, Pages 7913–7923, <https://doi.org/10.1093/nar/gky372>

Published: 10 May 2018 Article history ▾



Outils de prédition des points de branchement (BP) (1/5)



Pourquoi réaliser une comparaison des outils de prédition des BPs ?



6 outils sélectionnés : HSF, SVM-BPfinder, BPP, Branchpointer, LaBranchoR, RNABPS

¹Chow *et al.*, *Cell*, 1977

⁴Corvelo *et al.*, *PLOS Comput. Biol.*, 2010

⁷Signal *et al.*, *Bioinformatics*, 2018

²Shapiro *et al.*, *Nucleic Acids Res.*, 1987

⁵Mercer *et al.*, *Genome Res.*, 2015

⁸Paggi *et al.*, *RNA*, 2018

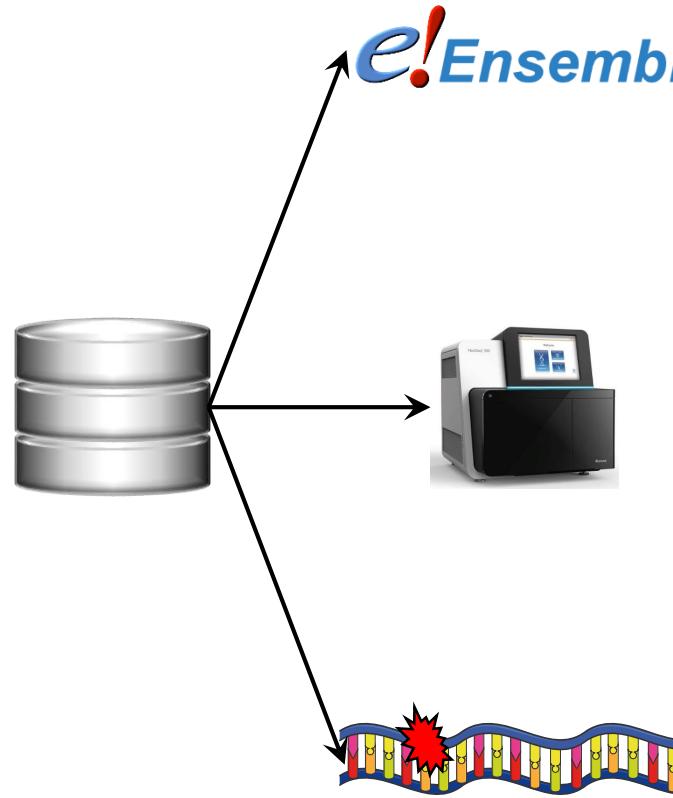
³Desmet *et al.*, *Nucleic Acids Res.*, 2009

⁶Zhang *et al.*, *Bioinformatics*, 2017

⁹Nazari *et al.*, *IEEE Access*, 2019



Jeux de données et performances (2/5)



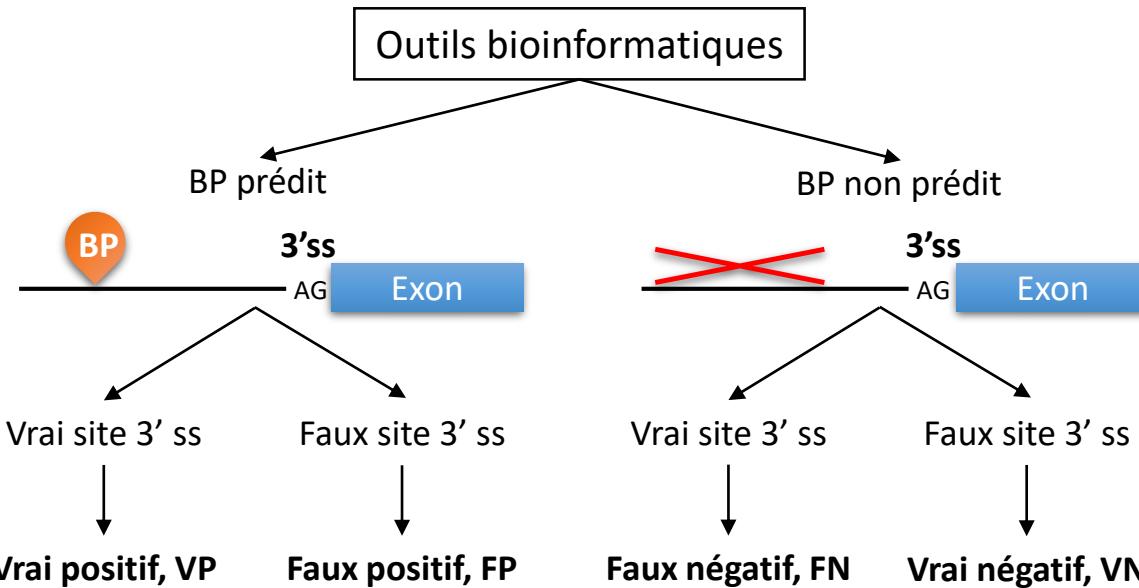
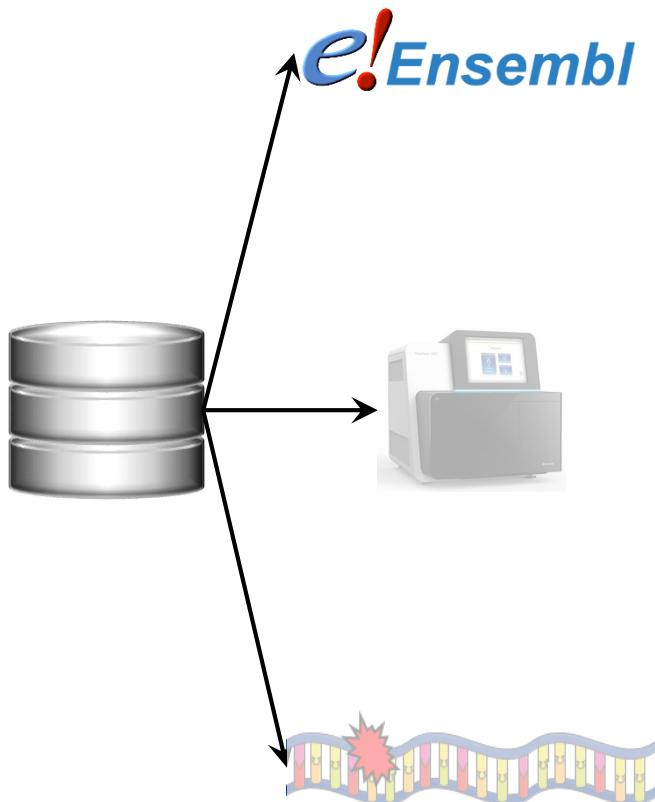
Sites physiologiques accepteurs
Sites négatifs tous motifs AG des 23 000 gènes

Sites accepteurs alternatifs par
RNA-seq → expression *versus*
prédition

Collection de variants avec étude ARN *in vitro*
dans la région BP (-44 à -18)



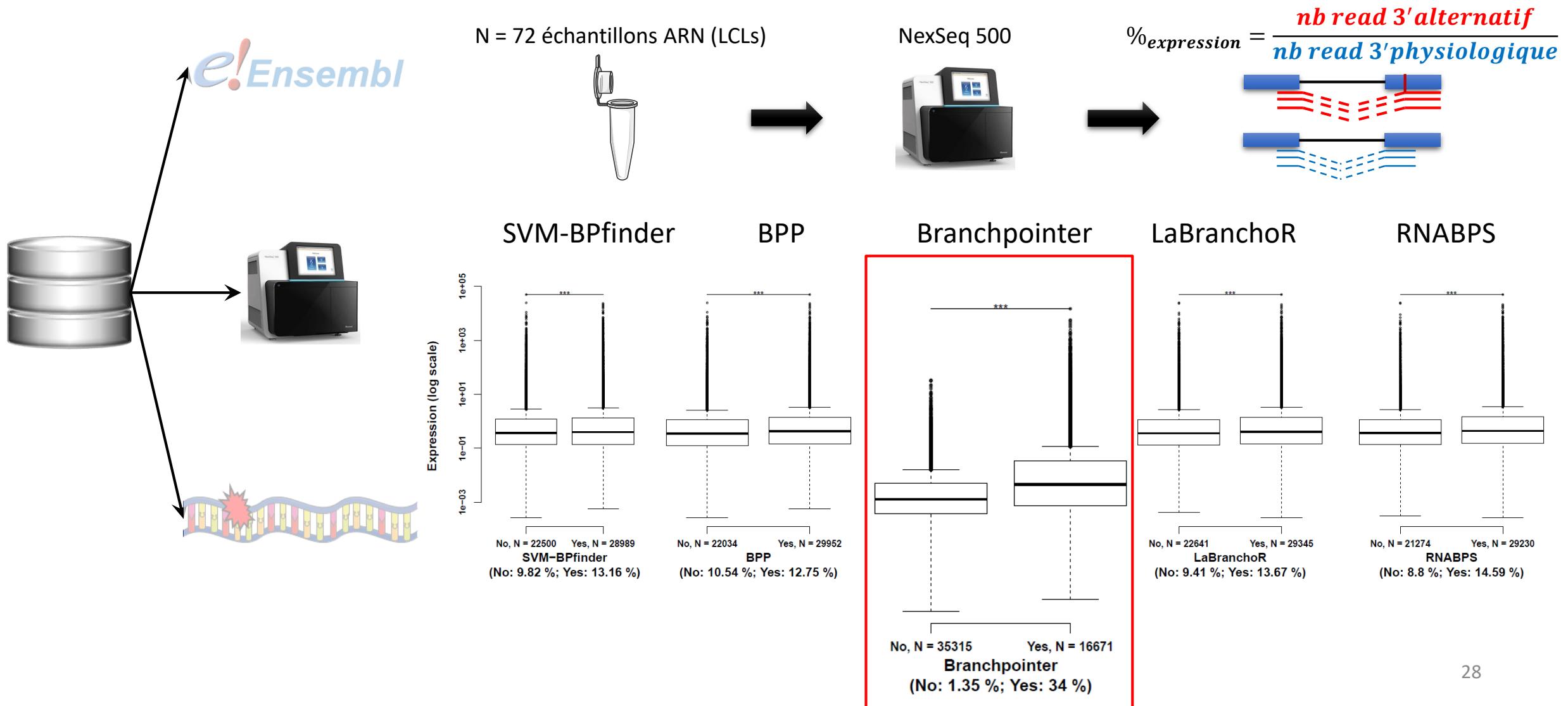
Jeux de données et performances (3/5)



	SVM-BPfinder	BPP	Branchpointer	LaBranchor	RNABPS
Exactitude	66.39 %	75.17 %	99.48 %	64.75 %	73.06 %
Sensibilité	66.39 %	75.17 %	93.27 %	64.75 %	73.05 %
Spécificité	66.39 %	75.17 %	99.49 %	64.75 %	73.06 %
VPP	0.45 %	0.70 %	30.06 %	0.43 %	0.62 %
VPN	99.88 %	99.92 %	99.98 %	99.87 %	99.91 %

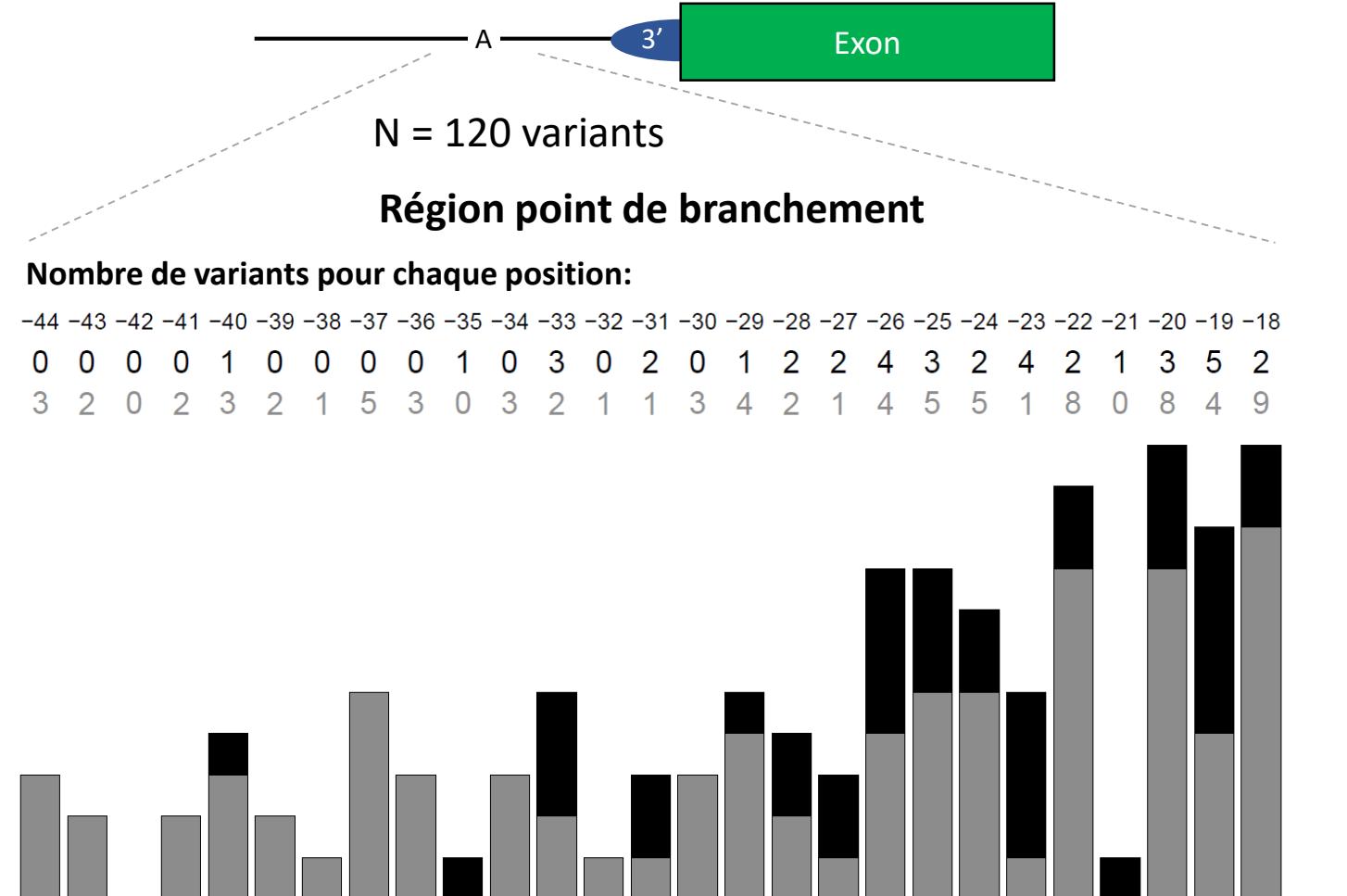
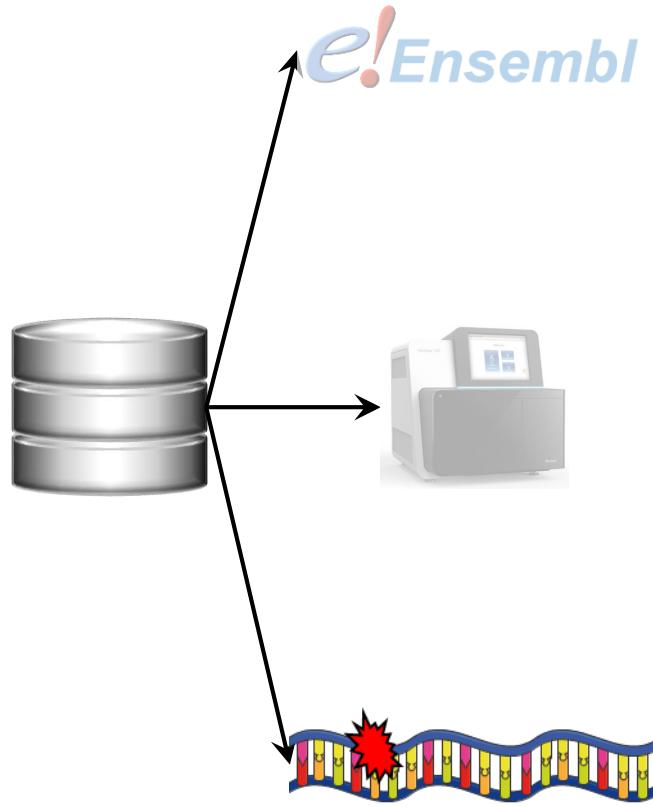


Jeux de données et performances (3/5)



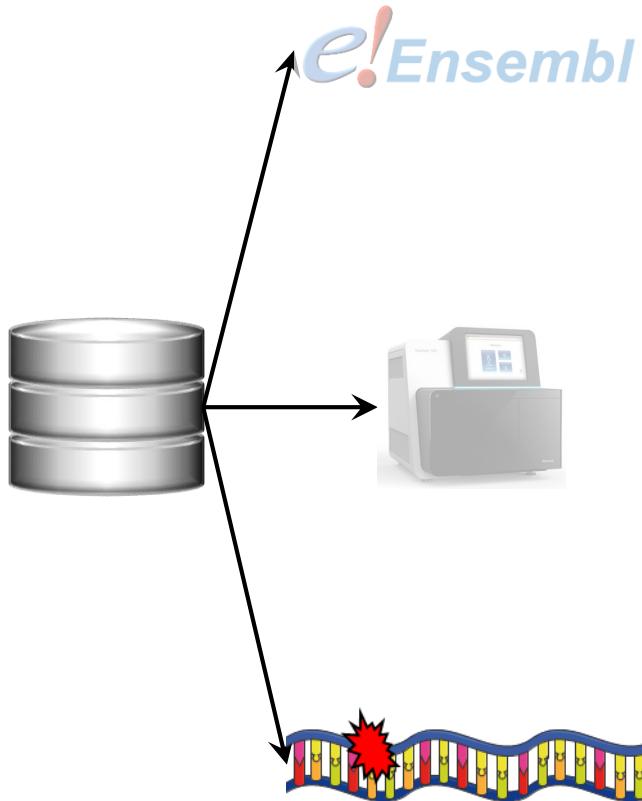


Jeux de données et performances (3/5)



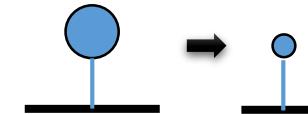


Jeux de données et performances (3/5)



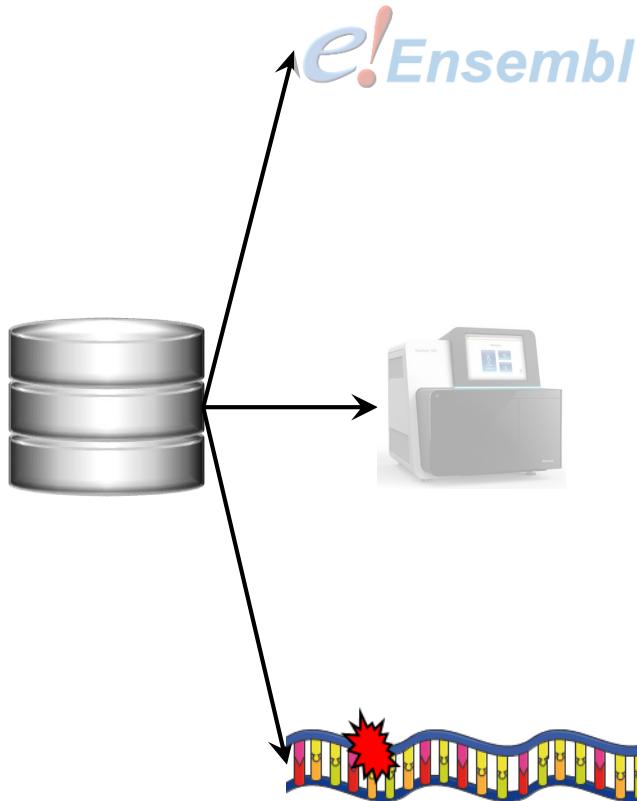
- Est-ce que le variant est situé dans un point de branchement ?
- Est-ce que le variant diminue le score du point de branchement ?

T R A Y *





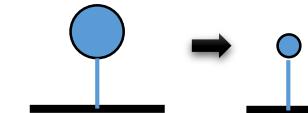
Jeux de données et performances (3/5)



- Est-ce que le variant est situé dans un point de branchement ?

T R A Y *

- Est-ce que le variant diminue le score du point de branchement ?



	SVM-BPfinder	BPP	Branchpointer	LaBranchoR	RNABPS
FP	6	7	12	15	12
FN	14	6	6	11	8
Exactitude	83.33 %	89.17 %	84.87 %	78.33 %	83.33 %
Sensibilité	63.16 %	84.21 %	84.21 %	71.05 %	78.95 %
Spécificité	92.68 %	91.46 %	85.19 %	81.71 %	85.37 %



Importance de l'apprentissage (4/5)

Détection d'un point de branchement parmi le bruit de fond

Branchpointer

A sequence of 'N' characters enclosed in brackets at both ends. A yellow circle with the letters 'BP' inside it highlights the second character from the left.



A sequence of 'N' characters enclosed in brackets at both ends, with a large red 'X' drawn across the entire sequence.

VS

Autres outils



A sequence of 'N' characters enclosed in brackets at both ends. A yellow circle with the letters 'BP' inside it highlights the second character from the left.

Altération d'un motif du point de branchement

BPP

Apprentissage sur les motifs les plus conservés sur plus de 200000 introns

VS

Autres outils

Apprentissage sur des BPs putatifs ou expérimentaux



Conclusion (5/5)

I. *Pouvoir utiliser les prédictions des points de branchement pour expliquer l'utilisation ou non d'un site accepteur*

- Meilleur score : **Branchpointer**

II. *Prédire l'effet d'un variant intronique sur un point de branchement*

- Meilleur score : **BPP**

RESEARCH ARTICLE *Epigenetics & Genomics*

Assessment of branch point prediction tools to predict physiological branch points and their alteration by variants

- › Raphael Leman, Hélène Tubeuf, Sabine Raad, Isabelle Tournier, Céline Derambure, Raphaël Lanos, Pascaline Gaildrat, Gaia Castelain, Julie Hauchard, Audrey Killian, Stéphanie Baert-Desurmont, Angelina Legros, Nicolas Goardon, Céline Quesnelle, Agathe Ricou, Laurent Castera, Dominique Vaur, Gérald Le Gac, Chandran Ka, Yann Fichou, Françoise Bonnet-Dorion, Nicolas Sevenet, Marine Guillaud-Bataille, Nadia Boutry-Kryza, Ines Schultz, Virginie Caux-Moncoutier, Maria Rossing, Logan C Walker, Amanda B Spurdle, Claude Houdayer, Alexandra Martins, Sophie Krieger

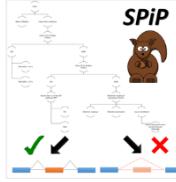
STATUS: IN REVISION

BMC Genomics

33



SPiP, Splicing Prediction Pipeline (1/4)



- **Objectif:**

Développer un outil de prédiction combinant les scores optimaux pour chaque motif

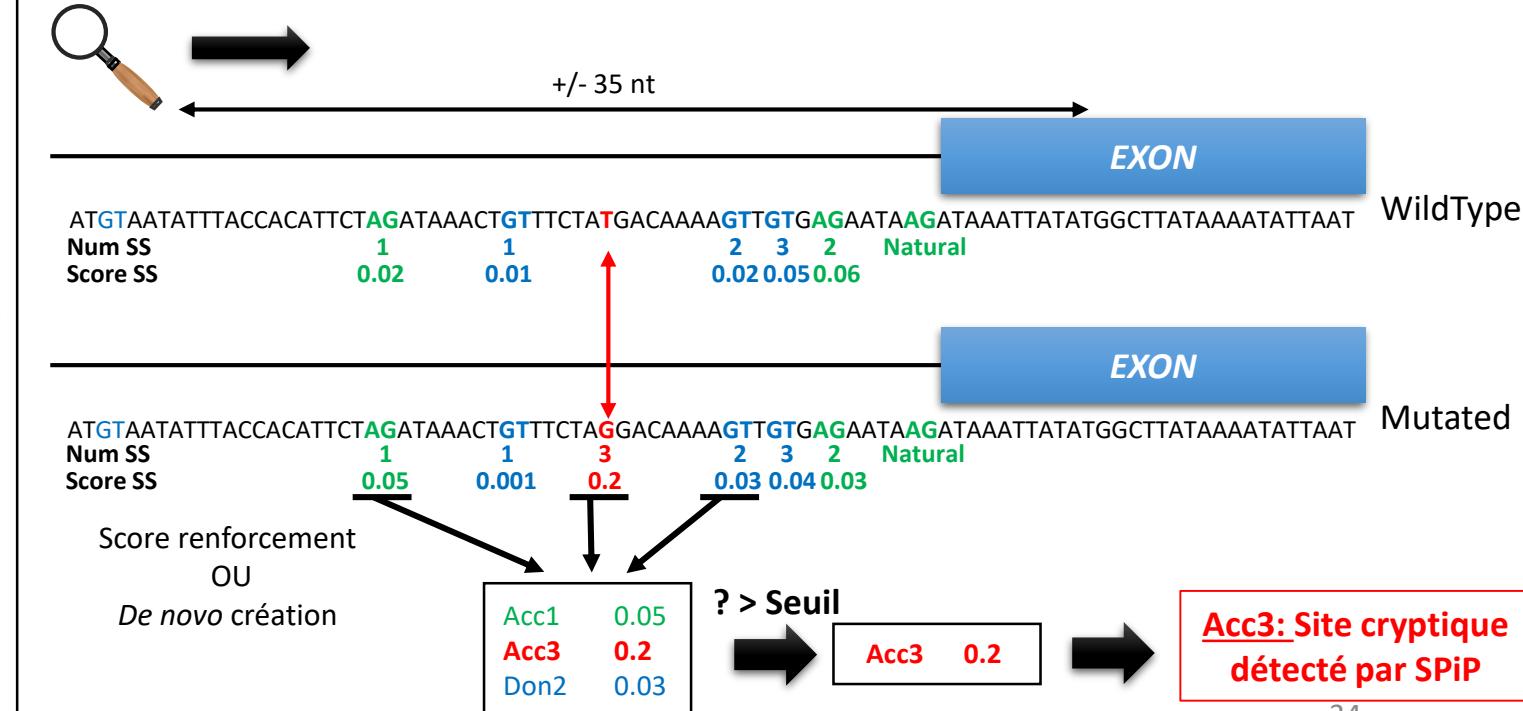
Comparaison des outils dans la littérature

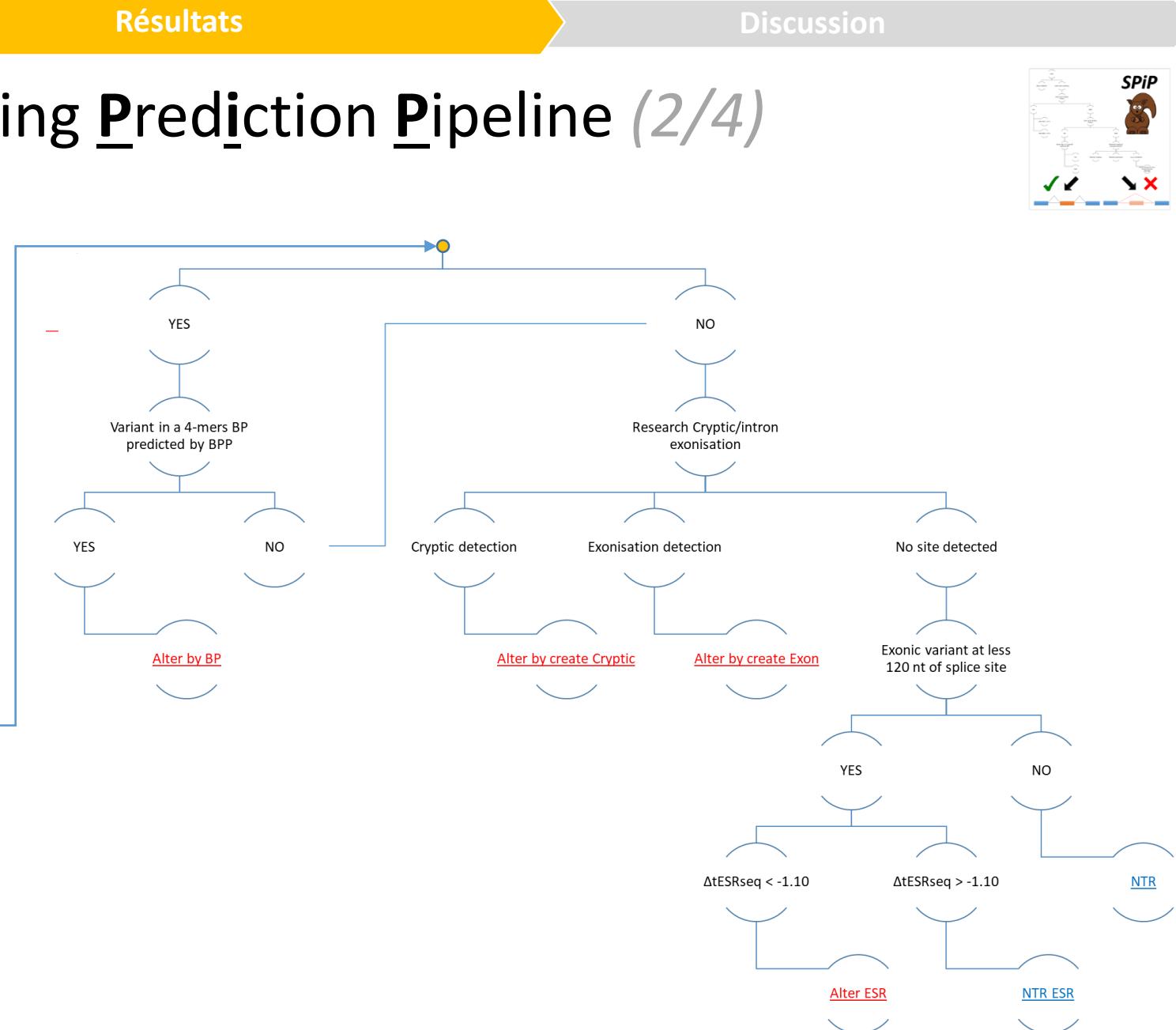
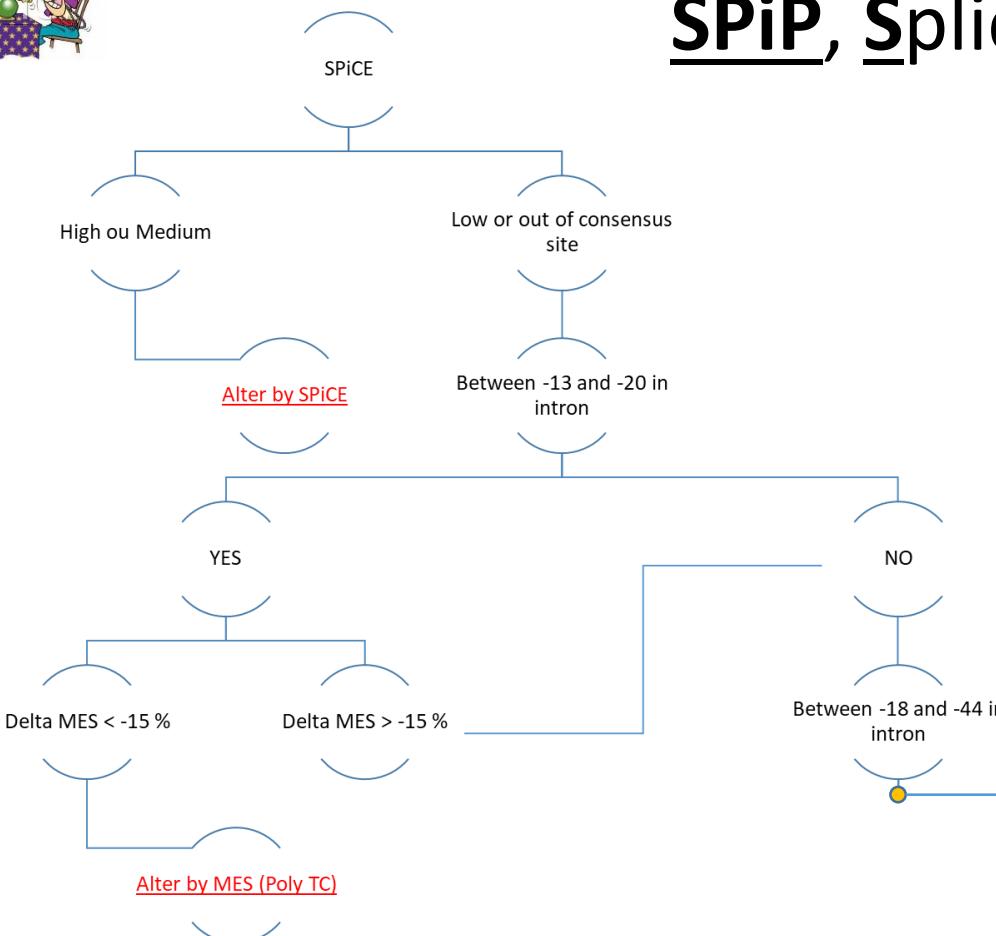
- ✓ **SPiCE** \leftrightarrow consensus sites¹
 - ✓ **BPP** \leftrightarrow branch points²
 - ✓ **ΔtESRseq** \leftrightarrow ESRs³
 - ✓ **MES** \leftrightarrow tract polypyrimidic⁴

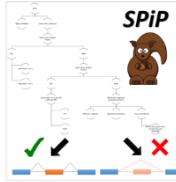
+

 e!Ensembl + Control AG/GT

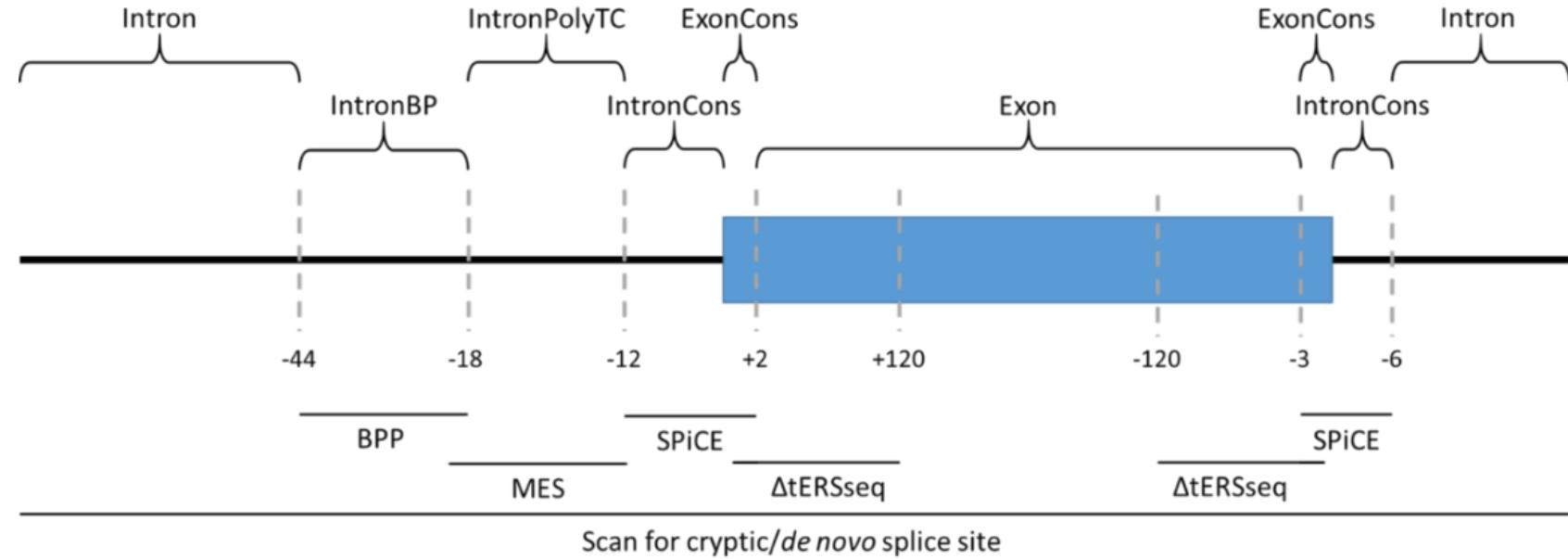
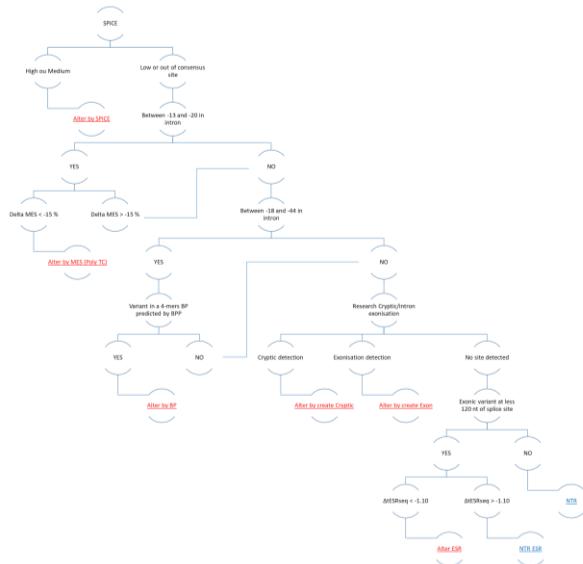
Nouvel outil pour la création d'un site d'épissage





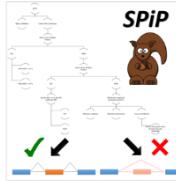


SPiP, Splicing Prediction Pipeline (2/4)

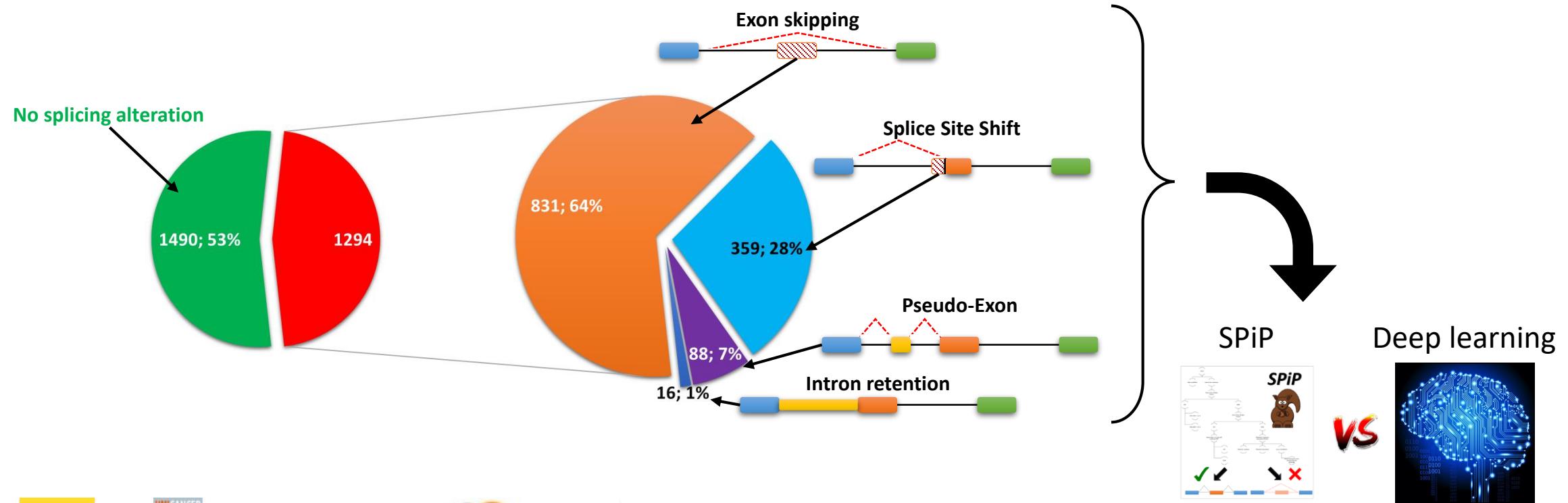




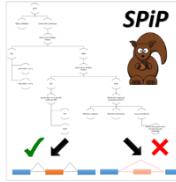
SPiP, Splicing Prediction Pipeline (3/4)



Jeux de 2 784 variants pour tester SPiP



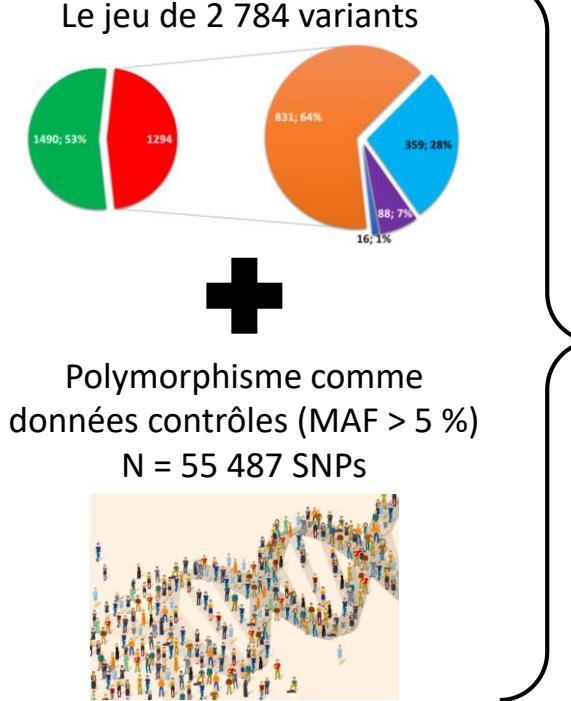
	EXACTITUDE	SENSIBILITÉ	SPÉCIFICITÉ
SPiP	80.21 %	90.96 %	70.87 %
SPANR ¹	75.45 %	78.37 %	72.32 %
SpliceAI ²	85.45 %	70.71 %	98.26 %



SPiP, Splicing Prediction Pipeline (4/4)

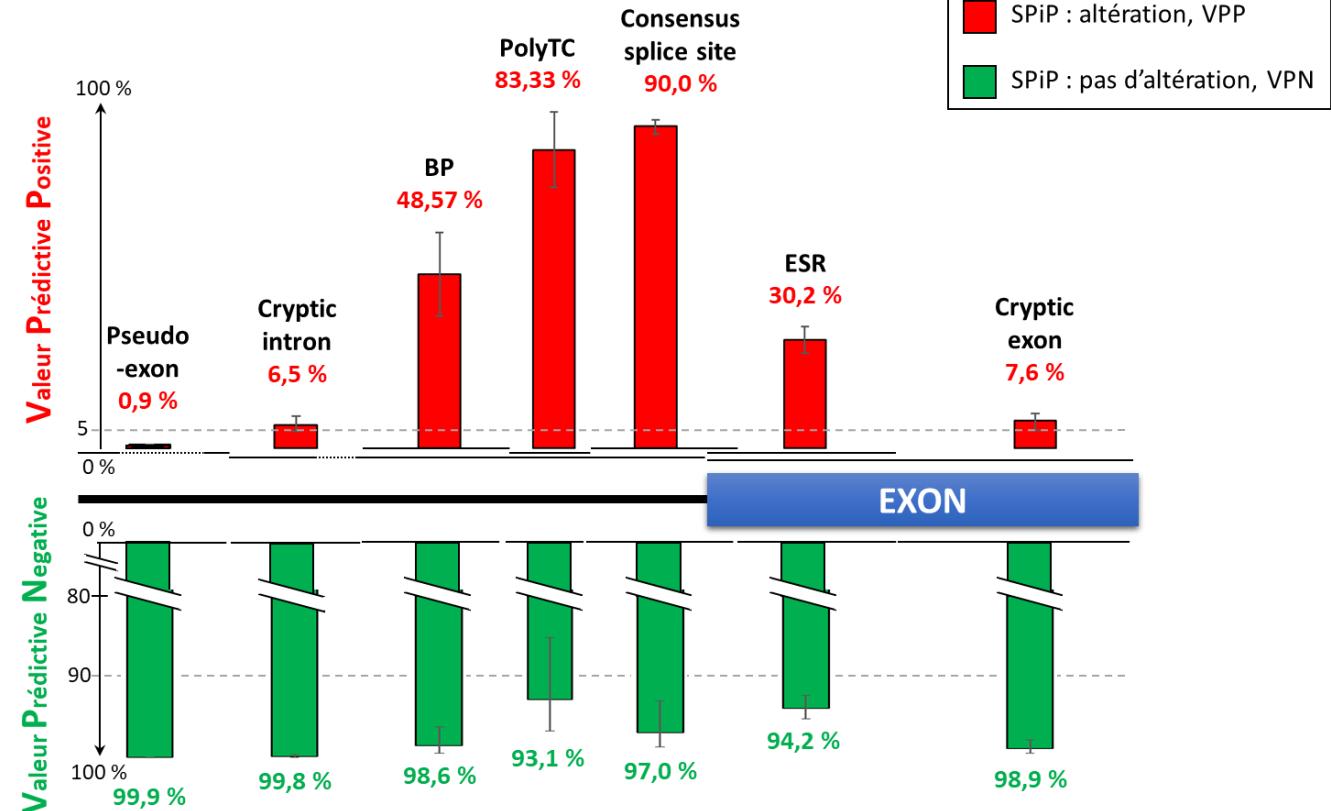
Estimation des probabilités d'altération de l'épissage selon :

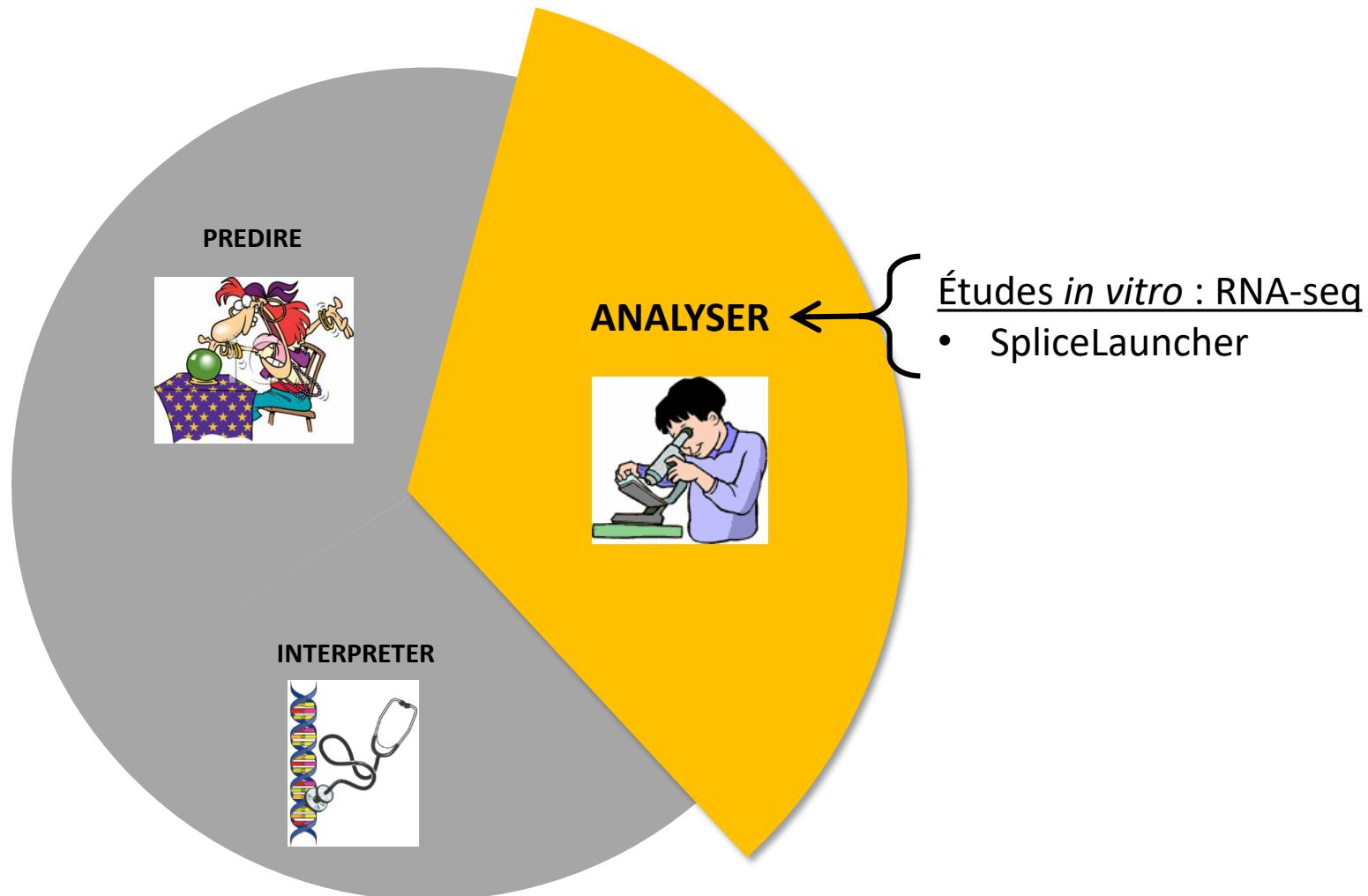
- 1) Position du variant
- 2) SPiP prédition



Théorème de Bayes

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

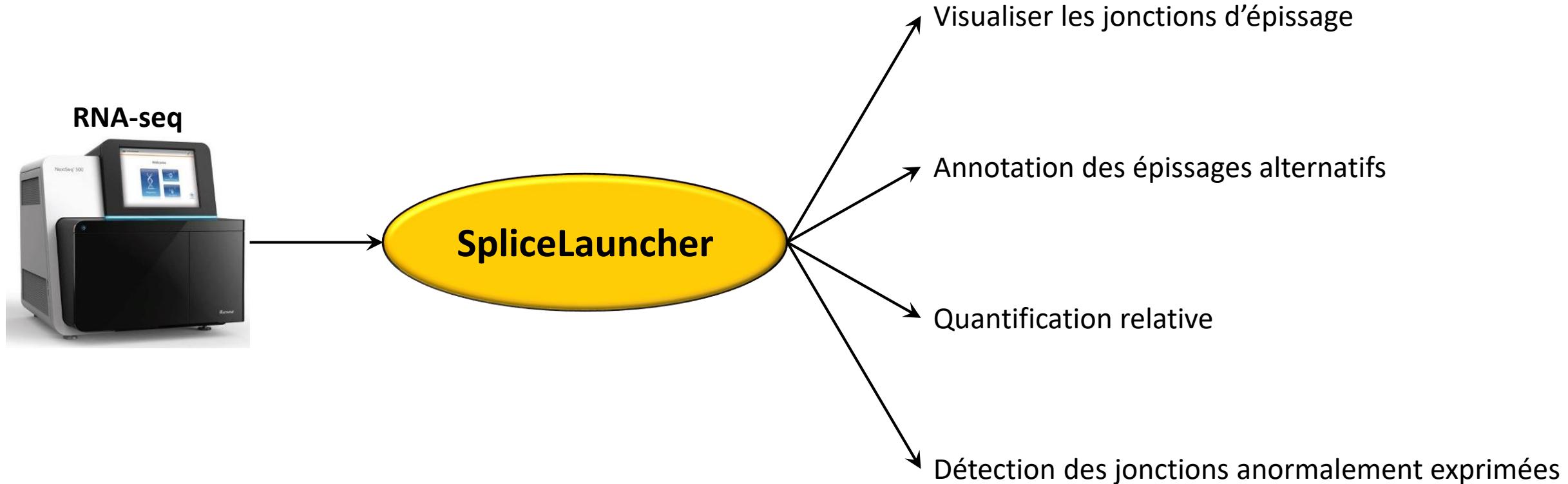






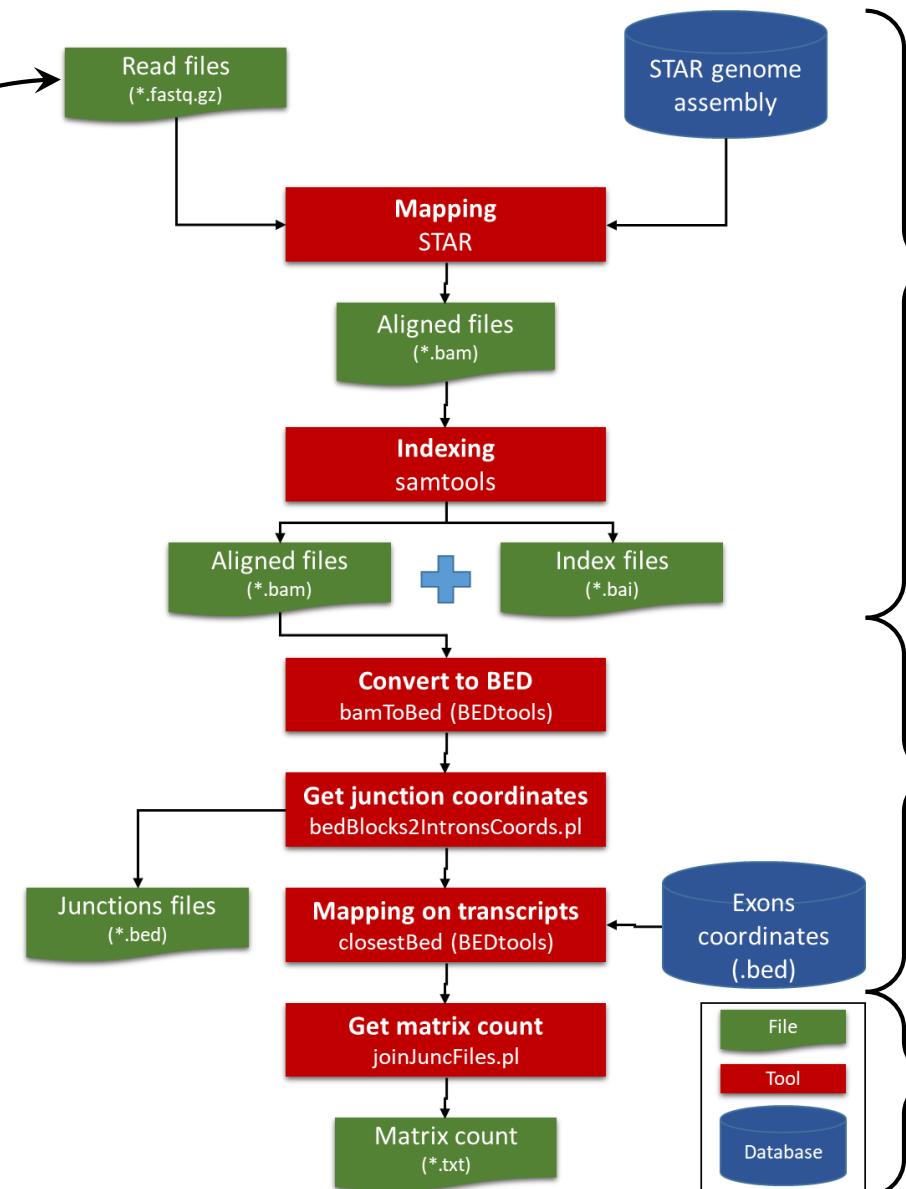
SpliceLauncher (1/4)

Plusieurs outils pour étudier l'expression des gènes et des exons **MAIS**
peu d'outils pour étudier la diversité des jonctions d'épissage

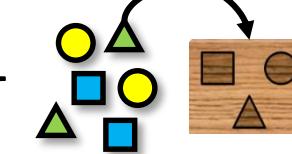




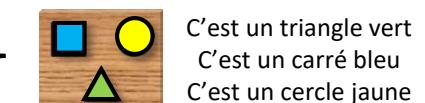
SpliceLauncher (2/4)



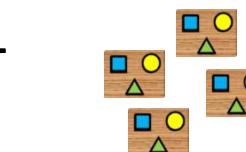
Alignment



Annotation



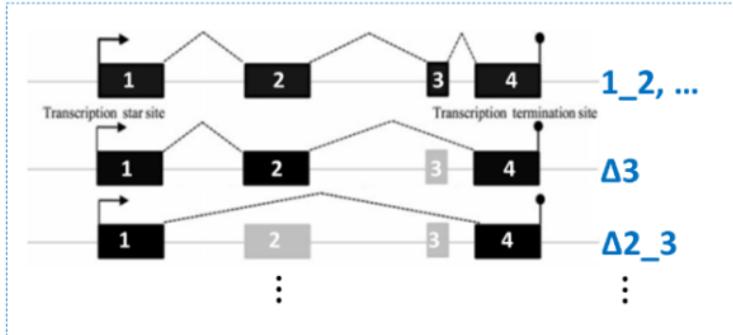
Comptage





SpliceLauncher (3/4)

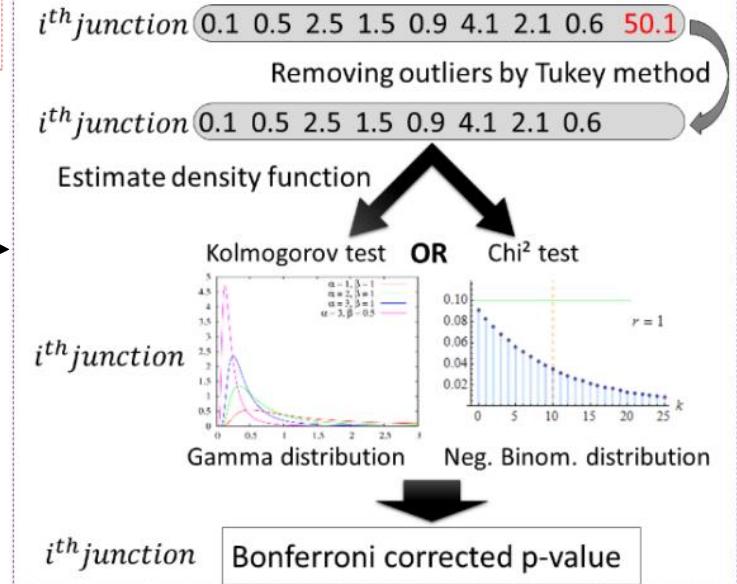
Annotation jonction



Expression relative

$E_x \quad E_{n-1} \quad E_n$	$\% = 100 \frac{(E_x - E_n)}{\frac{(E_x - E_{n-1}) + (E_{n-1} - E_n)}{2}}$
$E_x \quad E_{n-2} \quad E_{n-1} \quad E_n$	$\% = 100 \frac{(E_x - E_n)}{\frac{(E_x - E_{n-2}) + (E_{n-1} - E_n)}{2}}$
$E_x \quad I_a \quad I_b \quad E_n$	$\% = 100 \frac{(I_a - E_n)}{(E_x - E_n)}$
$E_x \quad I_a \quad I_b \quad E_n$	$\% = 100 \frac{(E_x - I_b)}{(E_x - E_n)}$

Analyse statistique





SpliceLauncher (4/4)

Excel file

Junction id	Gene	Read count	Relative expr.	Annotation
Chr1_123_456	ABC1	1 ... 5	0.1 ... 10.2	Δ3 c.23 c.230
⋮	⋮	⋮ ⋮ ⋮	⋮ ⋮ ⋮	⋮ ⋮ ⋮
ChrY_987_478	ZYZ2	10 ... 12	32.3 ... 25.4	Δ12p c.10 c.15



<https://github.com/raphaelleman/SpliceLauncher>



CORRECTED PROOF

SpliceLauncher: a tool for detection, annotation and relative quantification of alternative junctions from RNAseq data

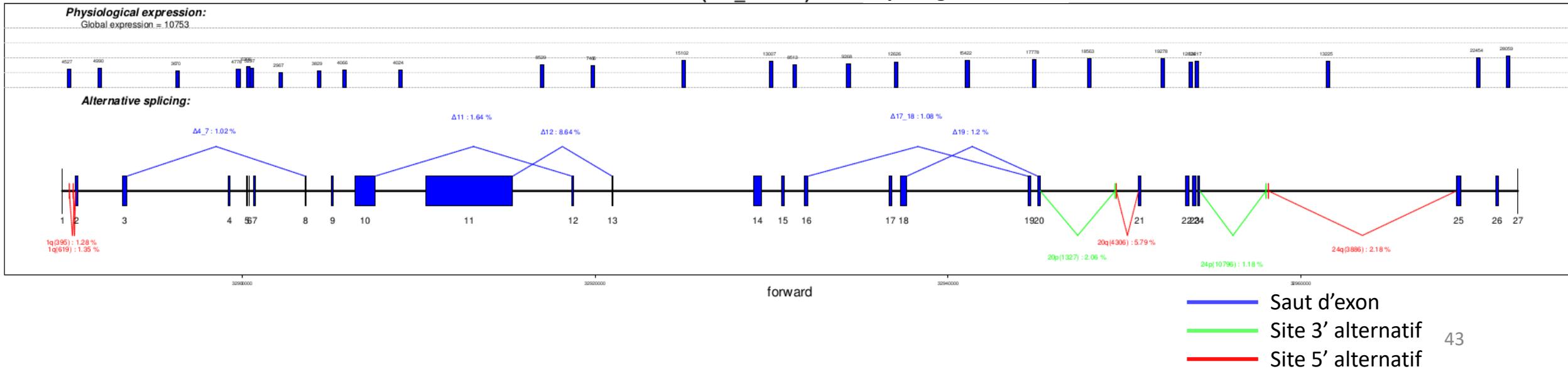
Raphaël Leman, Valentin Harter, Alexandre Atkinson, Grégoire Davy, Antoine Rousselin, Etienne Muller, Laurent Castéra, Frédéric Lemoine, Pierre de la Grange, Marine Guillaud-Bataille, Dominique Vaur, Sophie Krieger ✉

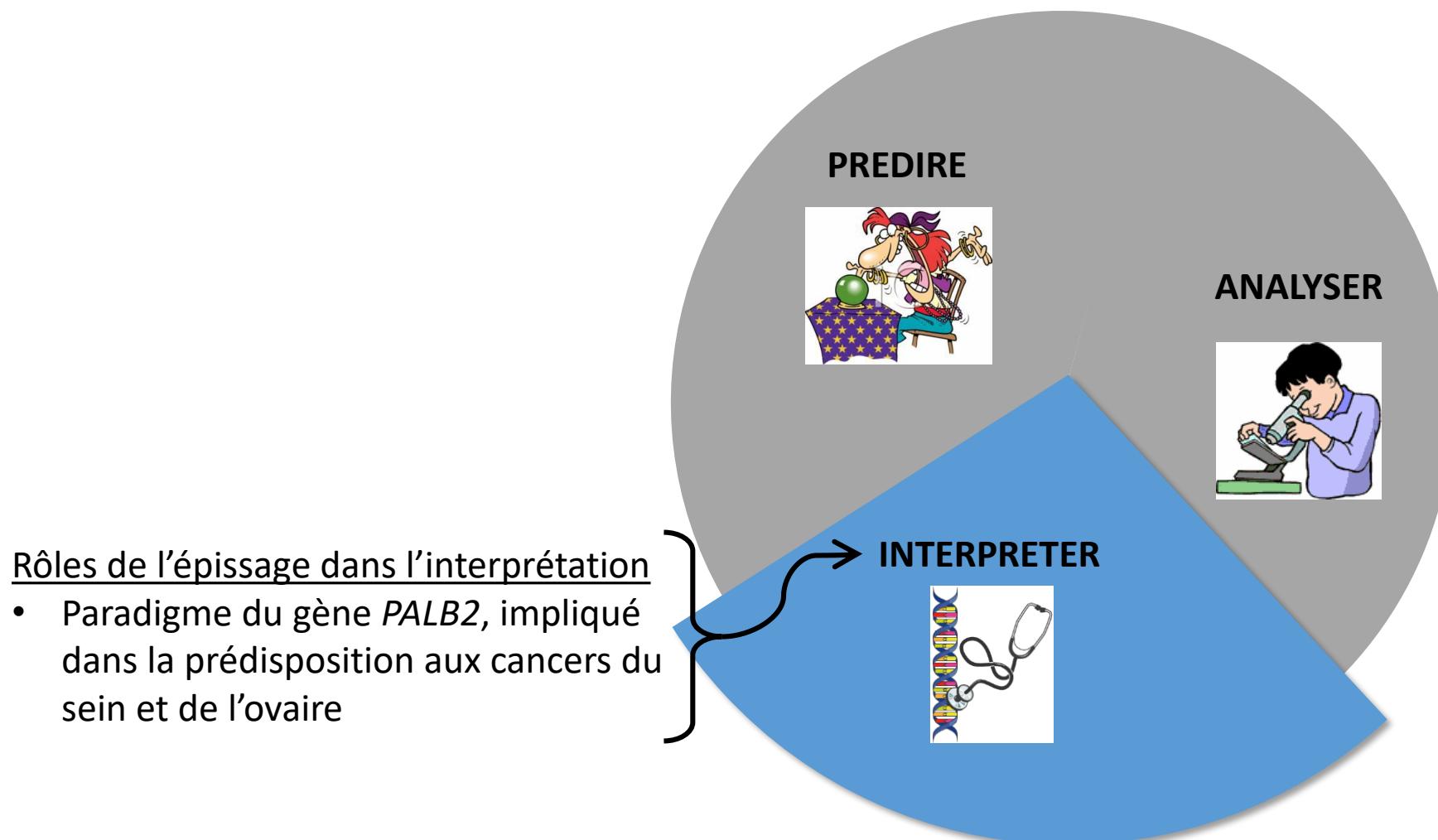
Bioinformatics, btz784, <https://doi.org/10.1093/bioinformatics/btz784>

Published: 16 October 2019 Article history ▾



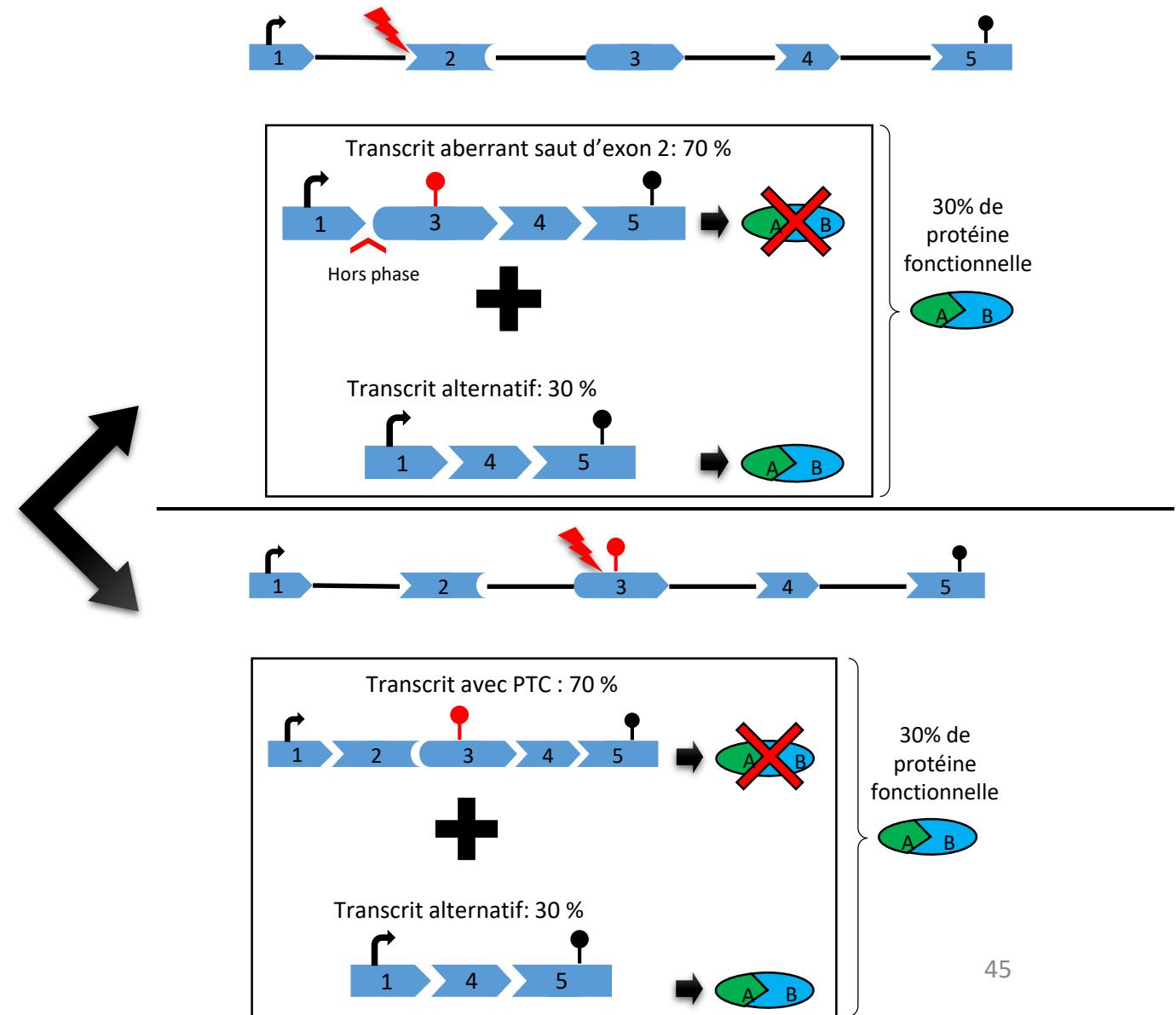
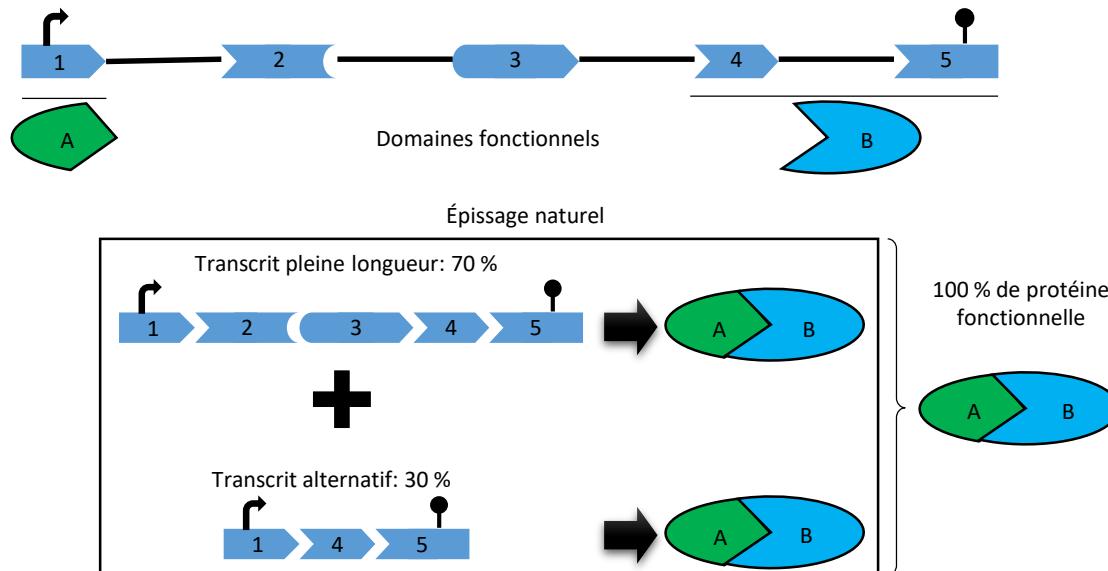
BRCA2 (NM_000059) > 1 % Epissage alternatif







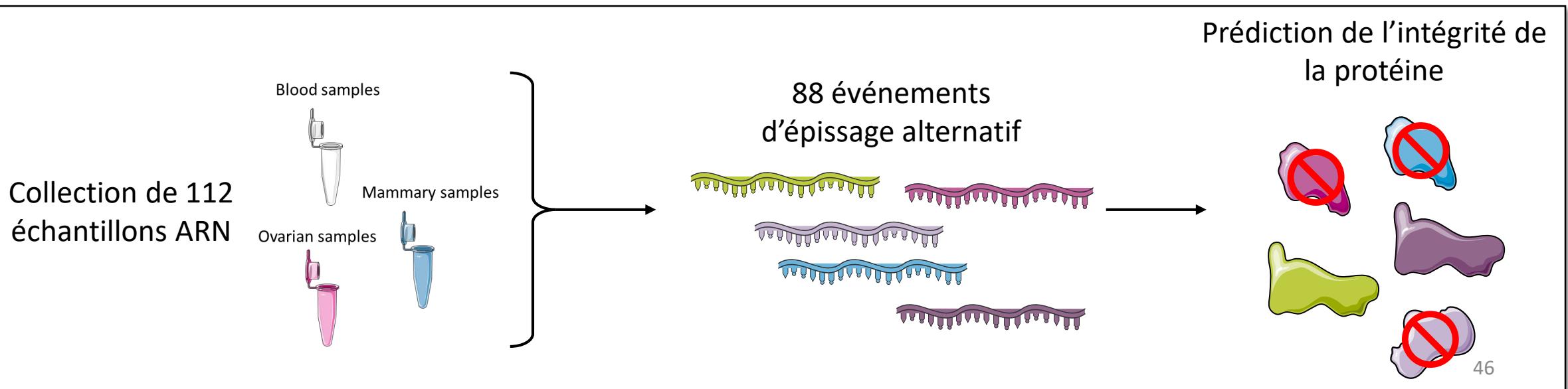
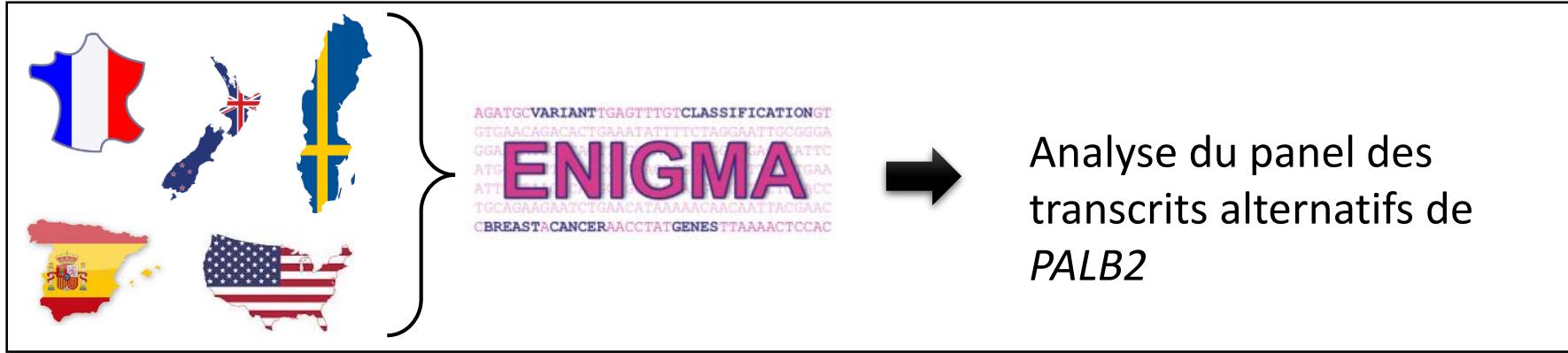
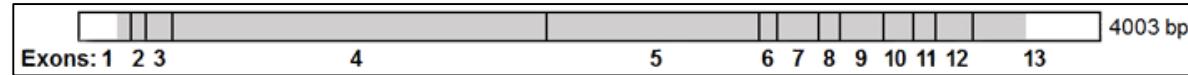
Épissage alternatif de *PALB2* et règles d'interprétation (1/3)





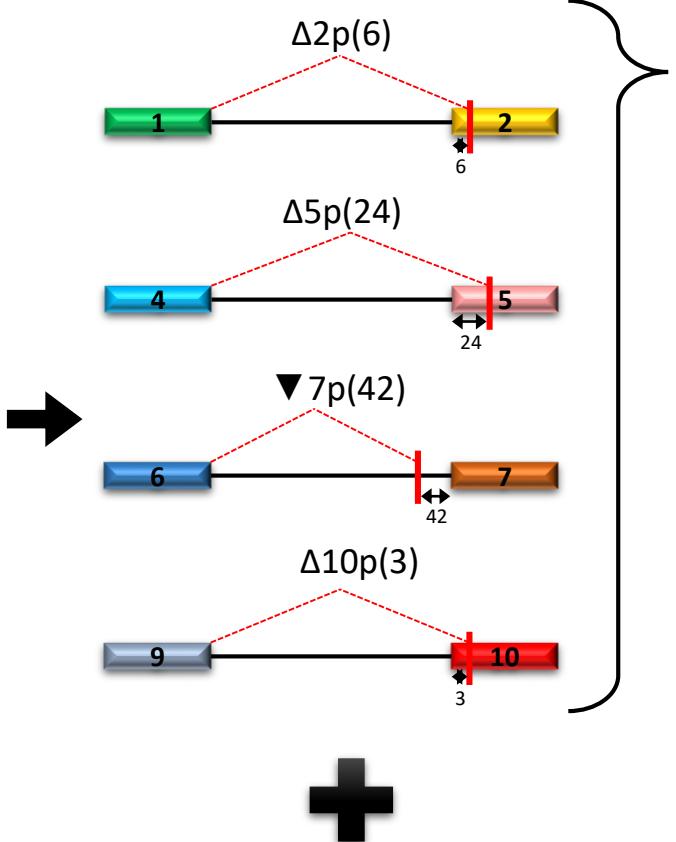
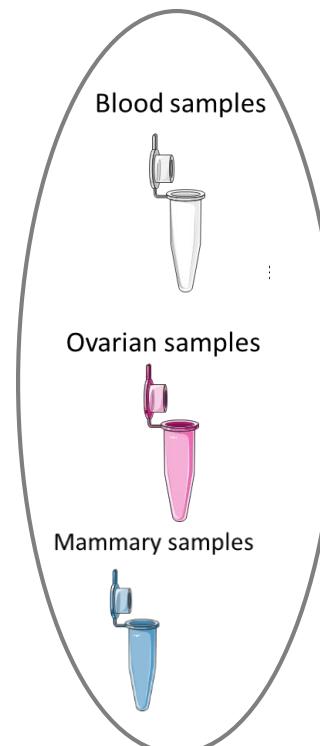
Épissage alternatif de *PALB2* et règles d'interprétation (2/3)

PALB2, gène de prédisposition aux cancers du sein et de l'ovaire





Épissage alternatif de *PALB2* et règles d'interprétation (3/3)



4 sites accepteurs pour lesquels la classe PVS n'est pas garantie



ACTGATGGTATGGGCCAAGAGATA
CAGGTACGGCTGTCACTTAGACCTCAC
CAGGGCTGGCATAAAGTCAGGGCAGAGC
CCATGGTGCATCTGACTCCTGAGGAGAA
GCAGGTTGGTACAAGTTACAAGACAGGT
GGCACTGACTCTCTGCCTATTGGTCTAT

ClinVar

c.49-2A>T
c.1685-2A>G
c.1685-1G>C
c.2587-2A>C
c.2997-2A>C

+
Pas d'épissage alternatif pouvant sauvegarder la protéine si variant tronquant

Journal of
Medical Genetics

Home / Archive / Volume 56, Issue 7

Latest content



Cancer genetics
Original article

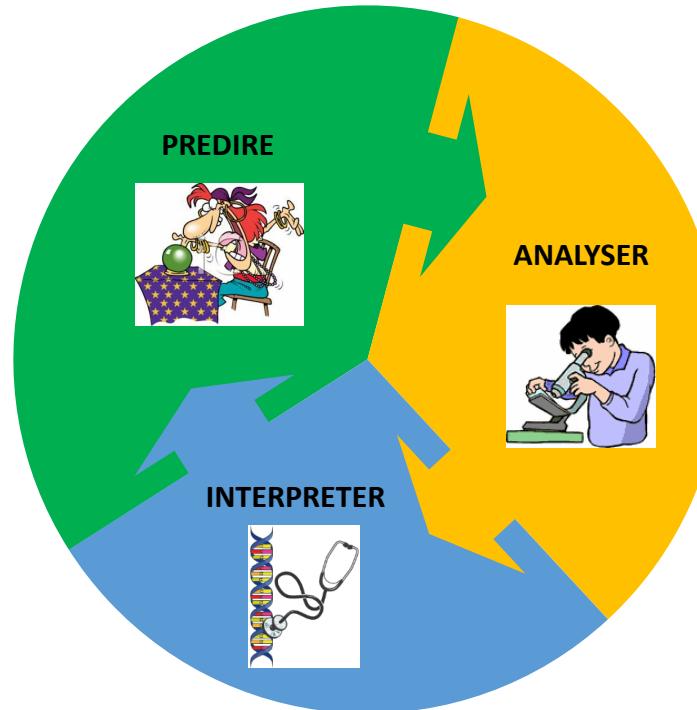
Alternative splicing and ACMG-AMP-2015-based classification of *PALB2* genetic variants: an ENIGMA report

Irene Lopez-Perolio¹, Raphaël Leman², Raquel Behar¹, Vanessa Lattimore³, John F Pearson³, Laurent Castéra², Alexandra Martins⁴, Dominique Vaur², Nicolas Goardon², Grégoire Davy², Pilar Garre¹, Vanesa García-Barberán¹, Patricia Llovet¹, Pedro Pérez-Segura¹, Eduardo Díaz-Rubio¹, Trinidad Caldés¹, Kathleen S Hruska⁵, Vickie Hsuan⁶, Sitao Wu⁶, Tina Pesaran⁶, Rachid Karam⁶, Johan Vallon-Christersson⁷, Ake Borg⁷, kConFab Investigators^{8,9}, Alberto Valenzuela-Palomo¹⁰, Eladio A Velasco¹⁰, Melissa Southee¹¹, Maaike P G Vreeswijk¹², Peter Devilee¹², Anders Kvist⁷, Amanda B Spurdle¹³, Logan C Walker³, Sophie Krieger², Miguel de la Hoya¹

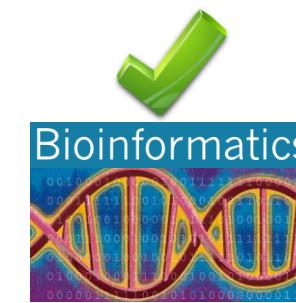
Résumé des travaux de thèse



SPiP → un outil pour prédire un défaut d'épissage quelle que soit la position du variant



Étude de l'épissage → une étape obligatoire pour l'interprétation des variants

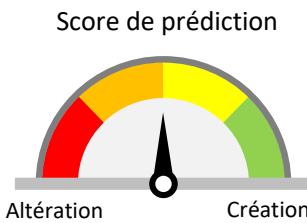


SpliceLauncher → Outils prometteurs pour analyser l'épissage



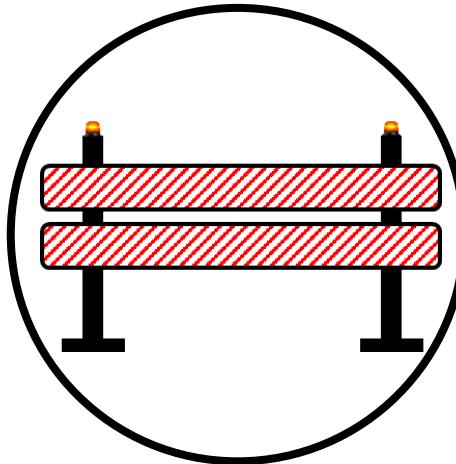
Limites actuelles des prédictions

Les prédictions d'épissage ne prédisent pas la pathogénicité

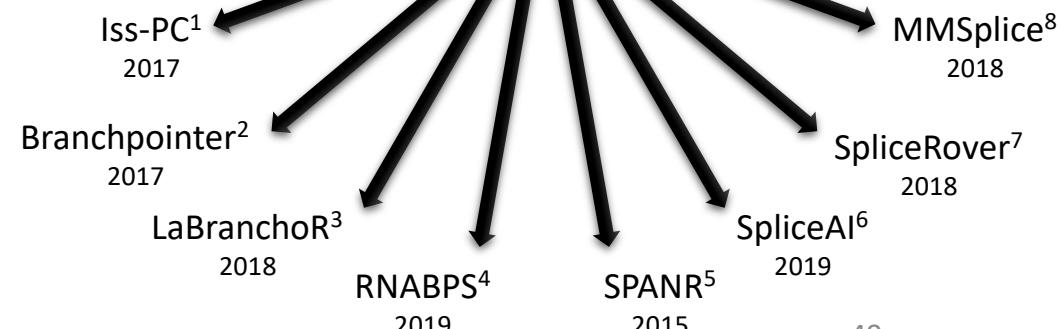
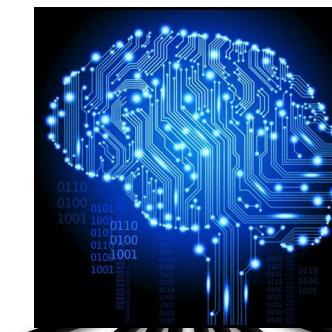
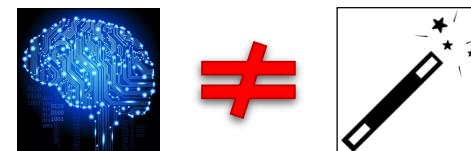


ACTGATGGTATGGGGCAAGAGATAATCT
CAGGTACGGCTGTCACTTAGACCTCAC
CAGGCTGGGCATAAAAGTCAGGGCAGAGC
CCATGGTGATCTGACTCTCTGAAGAGAAGT
GCAGGTTGGTATCAAGGTTACAAGACAGGT
GGCACTGACTCTCTGCCTATTGGTCTAT

ClinVar



Les limites à l'utilisation du *Deep Learning*



¹Xu et al., *Sci. Rep.*, 2017

⁵Xiong et al., *Science*, 2015

²Signal et al., *Bioinformatics*, 2017

⁶Jaganathan et al., *Cell*, 2019

³Paggi et al., *RNA*, 2018

⁷Zuallaert et al., *Bioinformatics*, 2018

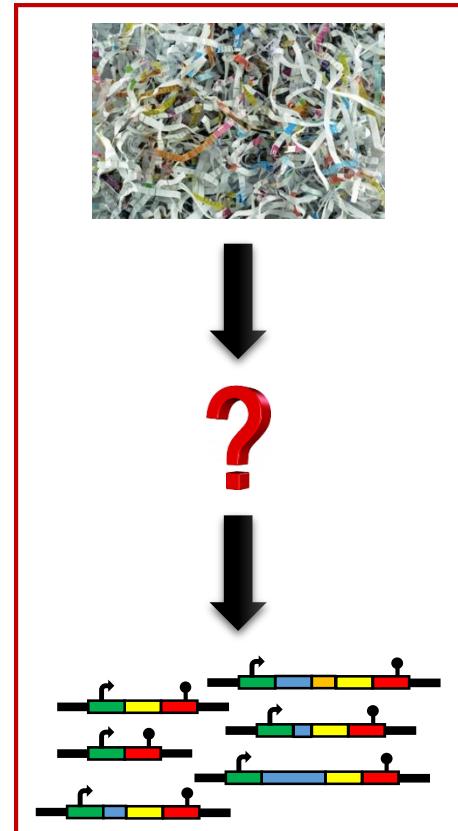
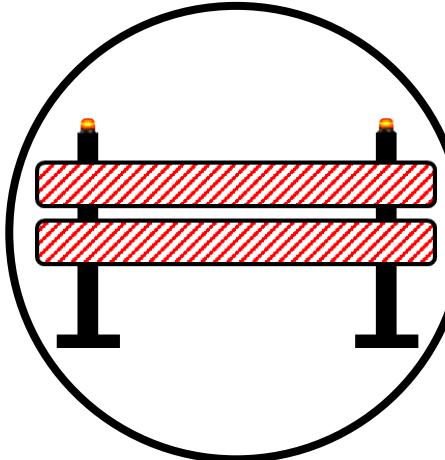
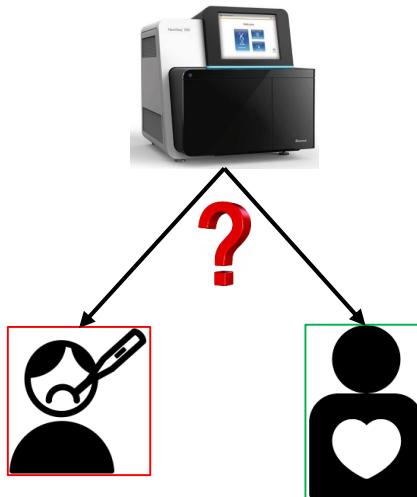
⁴Nazari et al., *IEEE Access*, 2019

⁸Cheng et al., *Genome Biol*, 2019

Limites actuelles des analyses par RNA-seq



Pas de recommandations pour
l'utilisation du RNA-seq en
diagnostic moléculaire

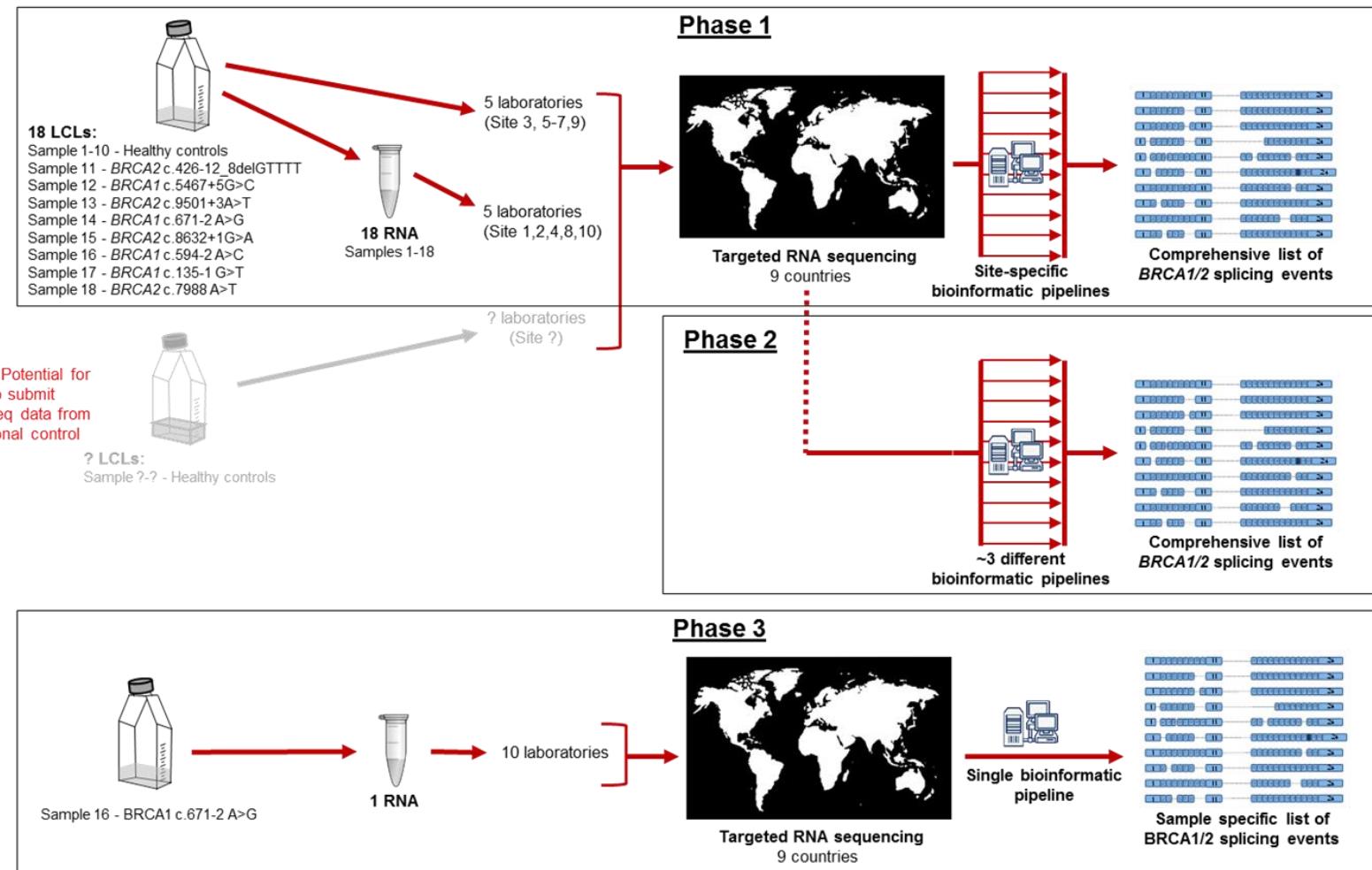




Projets en cours

– QC RNA-seq –

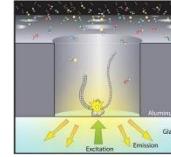
Établissement de recommandations pour l'analyse des défauts d'épissage par RNA-seq



ENIGMA

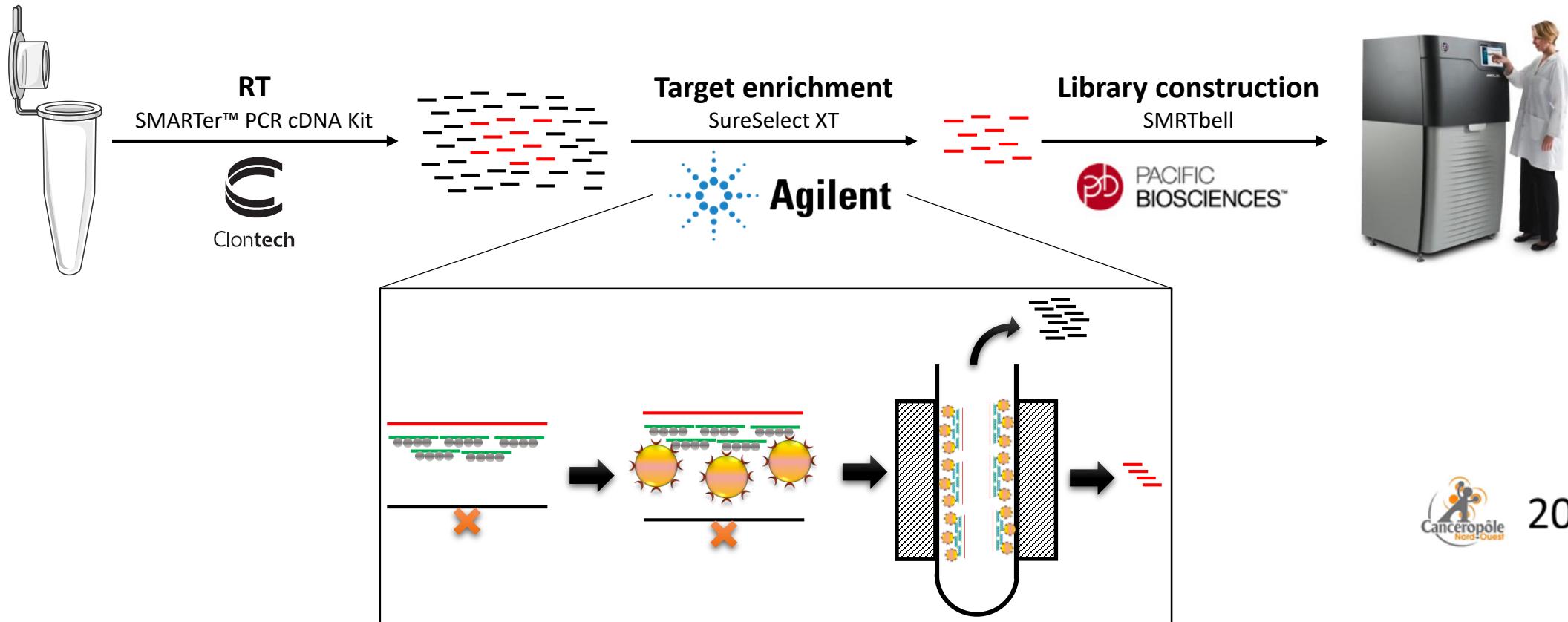
Projet mené par Logan Walker:

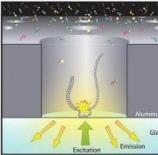
- 11 laboratoires
 - 9 pays
 - 9/11 ont fourni les données de séquençage
 - Données brutes analysées par SpliceLauncher



Projets en cours

– Séquençage long read pour RNA-seq ciblé –

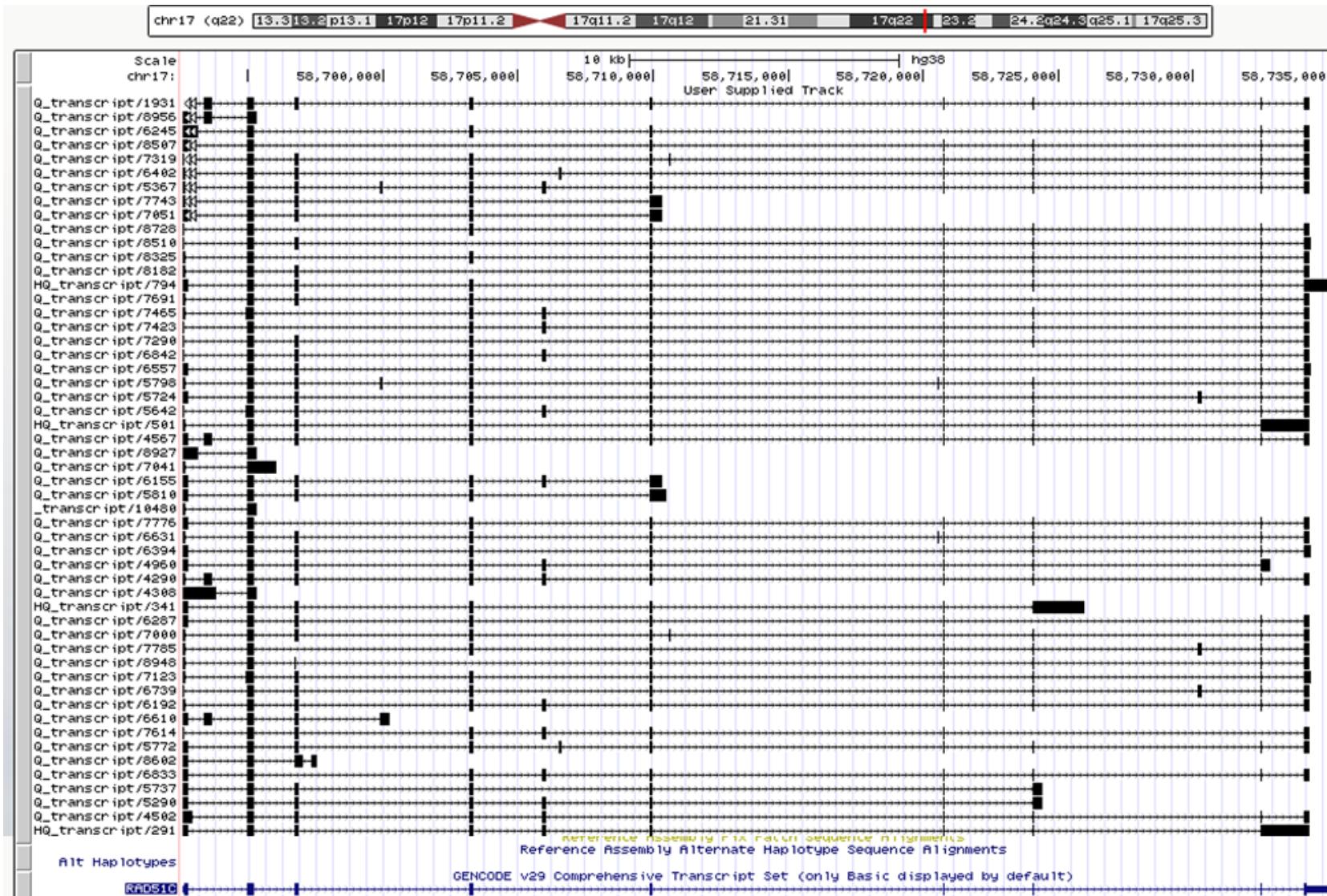




Projets en cours

– Séquençage long read pour RNA-seq ciblé –

RAD51C →

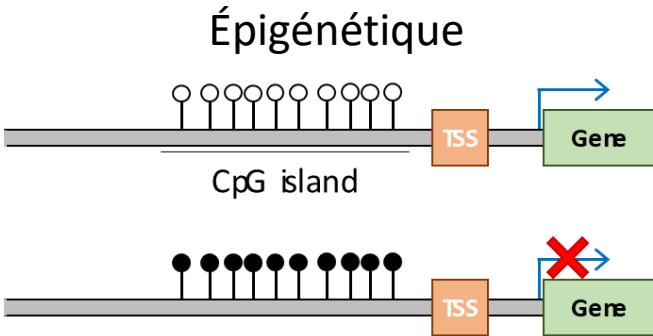
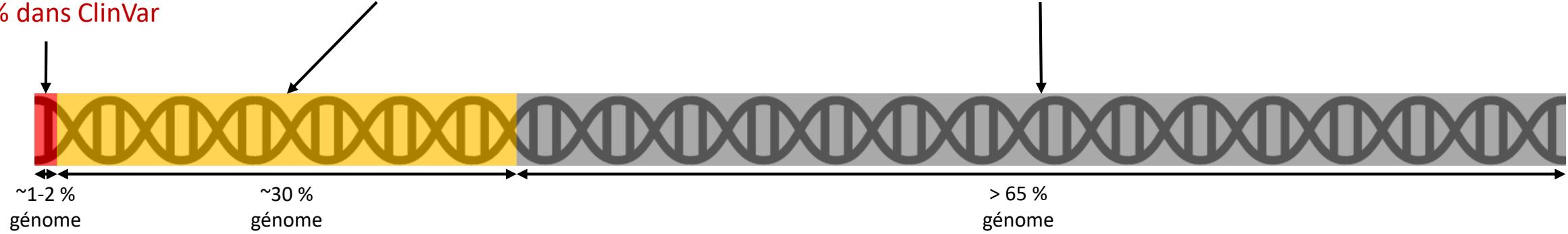


Conclusion et perspectives

Proportion de variants utilisables pour le diagnostic moléculaire

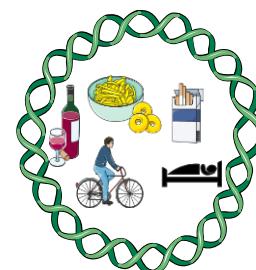
Régions exoniques

Variants non-sens
+ faux-sens
> 65 % dans ClinVar



Facteurs polygéniques

Gènes et environnement



Merci pour votre attention

~The End~

Sophie KRIEGER. Directrice de Thèse
Alexandra MARTINS. Co-directrice de Thèse

Inserm U1245

Dominique VAUR. Directeur du laboratoire de biologie, Caen
Claude HOUDAYER. Directeur du laboratoire de biologie, Rouen
Thierry FRÉBOURG. Directeur d'équipe, Inserm U1245

Agathe RICOU. Biogiste
Laurent CASTERA. Biogiste
Stéphanie BAERT-DESURMONT. Biogiste
Etienne MULLER. Assistant biologiste
Isabelle TOURNIER. Maître de conférences
Angelina LEGROS. Chef d'équipe
Nicolas GOARDON. Cadre de service
Omar SOUKARIEH. Ingénieur
Pascaline GAILDRAT. Chargée de recherche
Alexandre ATKINSON. Bioinformaticien

Antoine ROUSSELIN. Bioinformaticien
Raphaël LANOS. Bioinformaticien
Céline DERAMBURE. Technicienne
Gaia CASTELAIN. Technicienne
Julie HAUCHARD. Technicienne
Audrey KILLIAN. Technicienne
Céline QUESNELL. Technicienne
Grégoire DAVY. Etudiant en thèse
Hélène TUBEUF. Etudiant en thèse
Sabine RAAD. Etudiant en thèse

Inserm UMR1078

Claude FEREZ. Directeur d'équipe
Chandran KA. Ingénieur
Gérald Le GAC. Maître de conférences
Yann FICHOU. Ingénieur
Marie-Pierre AUDREZET. Biogiste

Inserm U1016

Béatrice PARFAIT. Biogiste
Dominique VIDAUD. Biogiste
Emmanuelle GIRODON. Biogiste
Laurence PACOT. Interne

Valentin HARTER. Biostatisticien
Pierre de la GRANGE. GenoSplice
Florence RIANT. Biogiste
Frédéric LEMOINE. Bioinformaticien

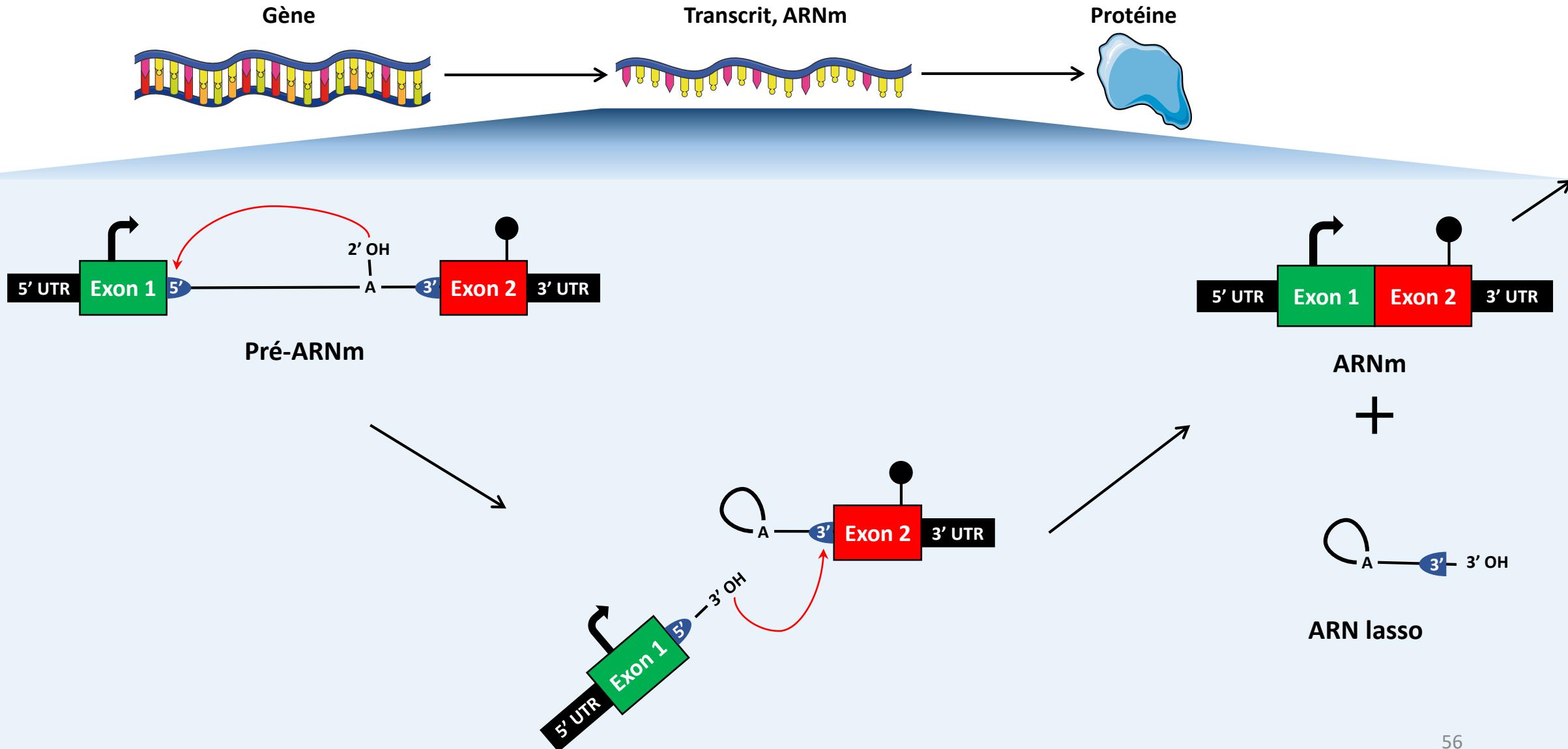
Réseau GGC

Brigitte BRESSAC-DE PAILLERETS. Biogiste
Capucine DELNATTE. Biogiste
Danielle MULLER. Biogiste
Etienne ROULEAU. Biogiste
Florence COULET. Biogiste
Virginie CAUX-MONCOUTIER. Ingénieur
Violaine BOURDON. Biogiste
Sylvie MAZOYER. Biogiste
Sandrine M. CAPUTO. Chargée de recherche
Nicolas SEVENET. Biogiste
Françoise BONNET-DORION. Ingénieur
Mélanie LÉONE. Biogiste
Marine GUILLAUD-BATAILLE. Ingénieur
Françoise REVILLON. Biogiste
Nadia BOUTRY-KRYZA. Biogiste
Joanna SOKOLOWSKA. Biogiste
Inès SCHULTZ. Biogiste

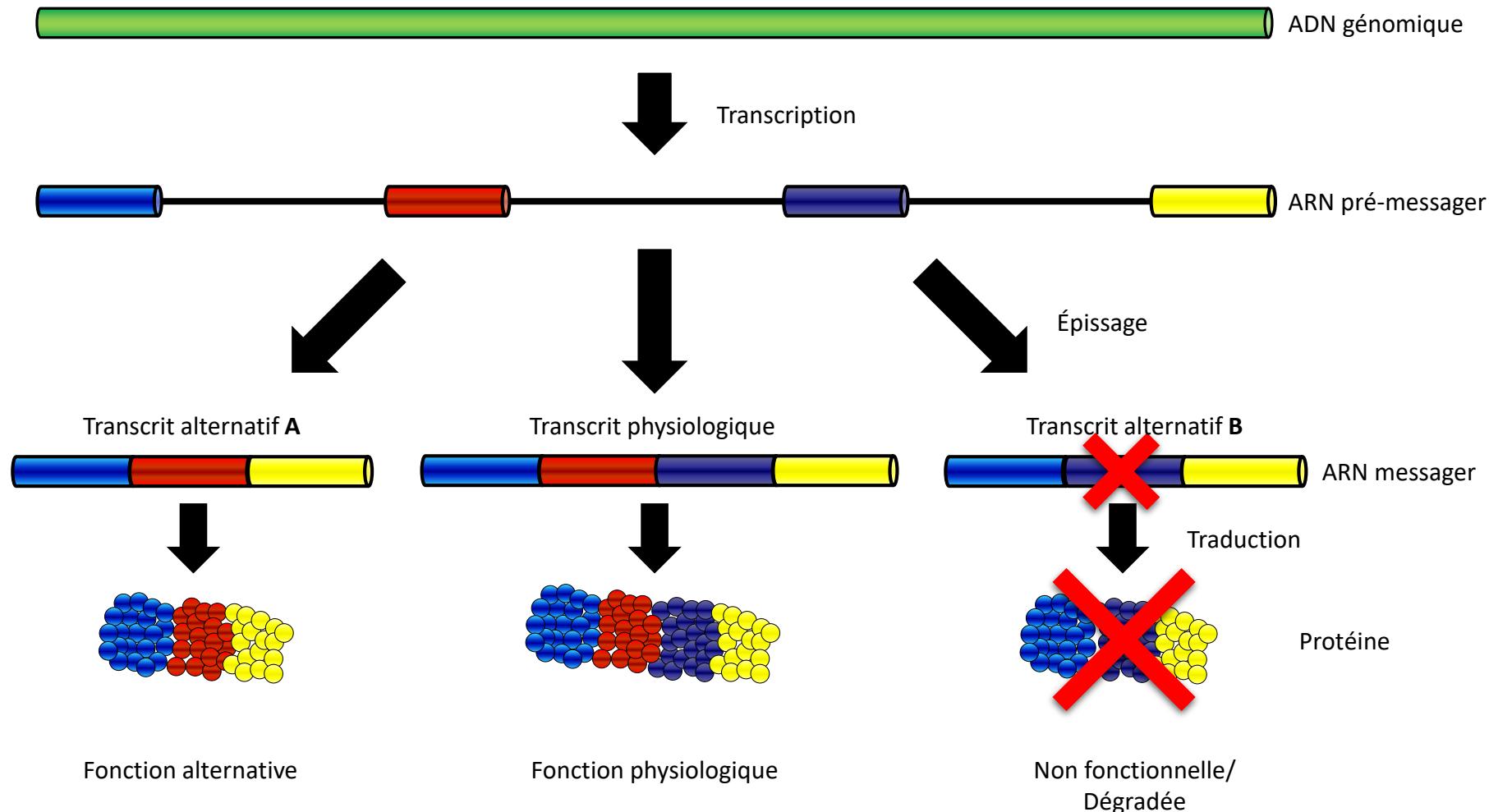
ENIGMA

Amanda B. SPURDLE. Professor
Anders KVIST. Researcher
Barbara WAPPENSCHMIDT. Fachhumangenetikerin
Irene LOPEZ-PEROLIO. Research Associate
Logan C WALKER. Associate Professor
Maaike P G VREESWIJK. Researcher
Maria ROSSING. Researcher
Melissa SOUTHEY. Professor
Michael T. PARSONS. Research assistant
Miguel de la HOYA. Biogiste
Rachid KARAM. Ambry genetics
Vanessa LATTIMORE. Research Fellow

Mécanisme d'épissage



L'épissage alternatif

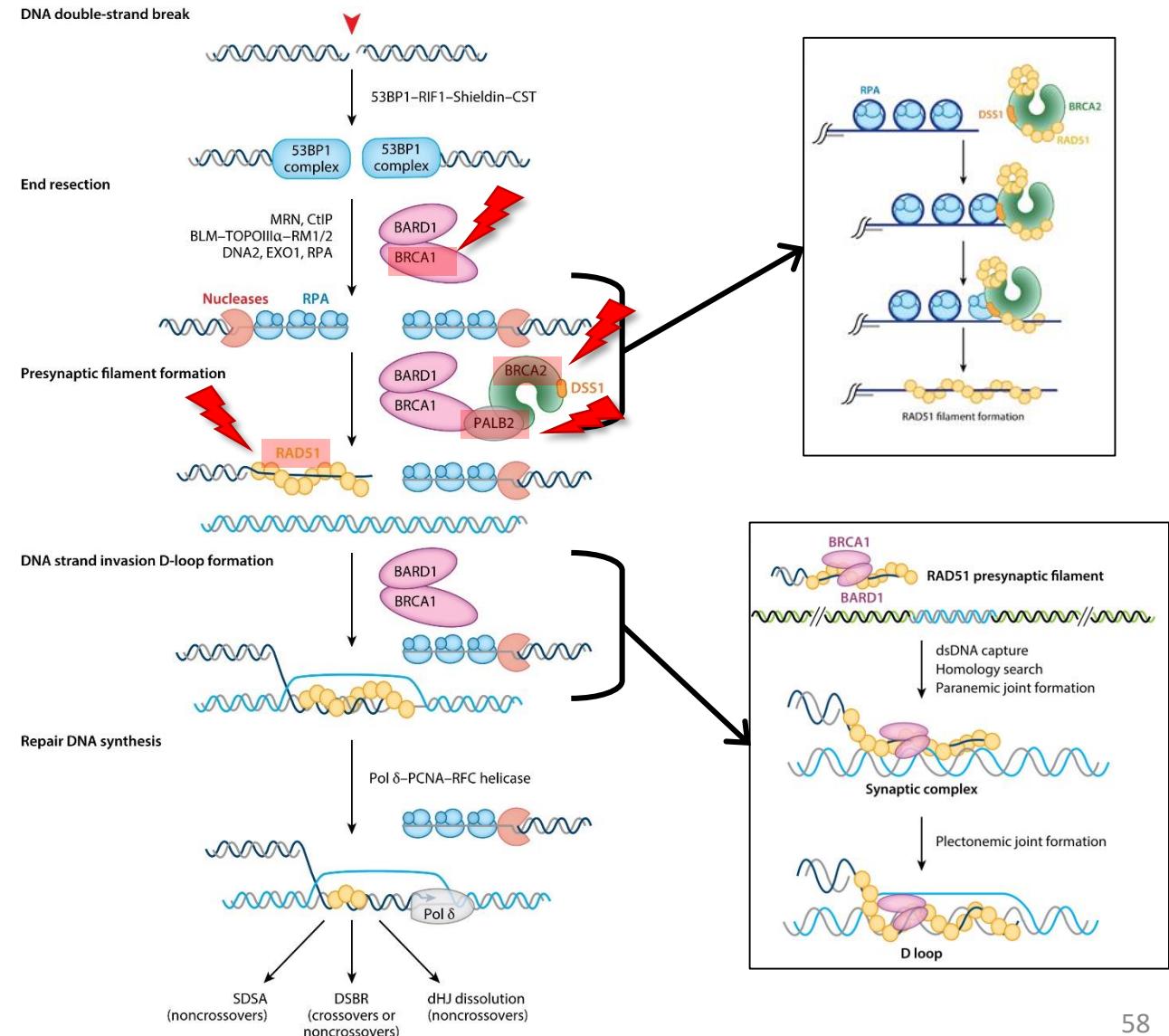


Gènes impliqués dans le syndrome HBOC

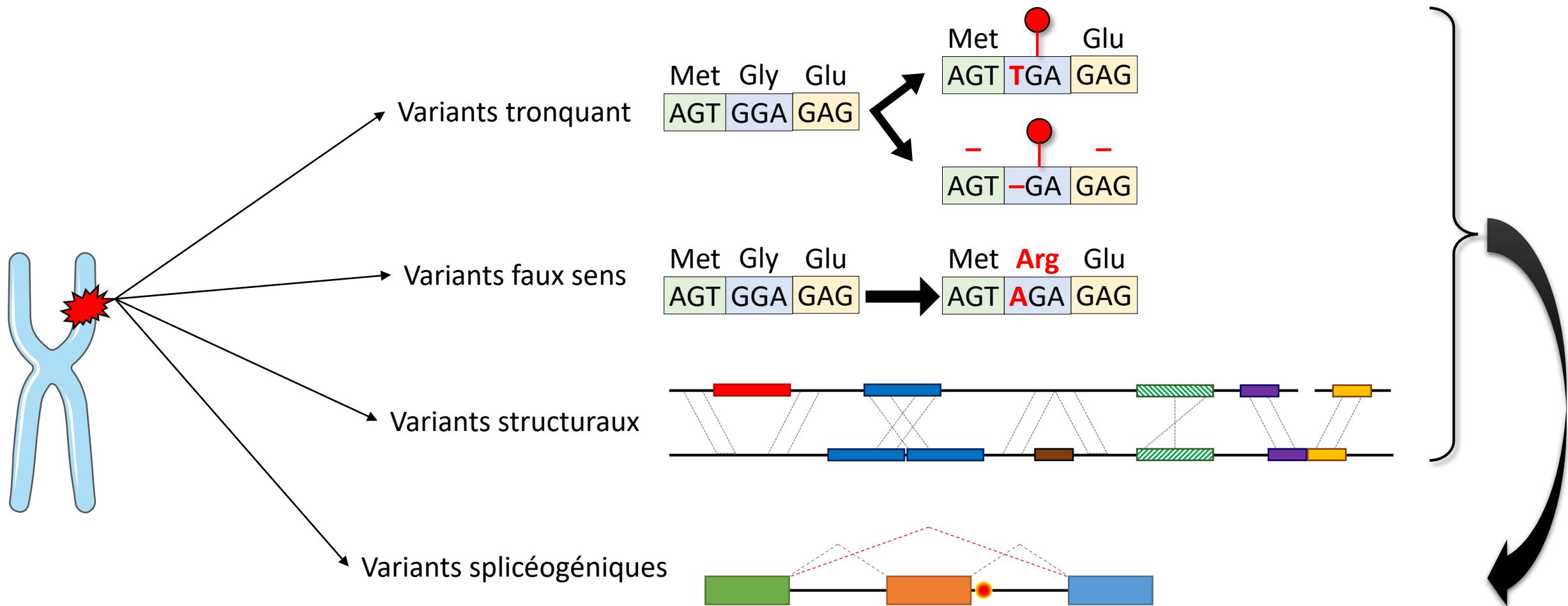
5 genes majeurs de predisposition :

- *BRCA1* (1994)
- *BRCA2* (1995)
- *PALB2* (2007)
- *RAD51C* (2010)
- *RAD51D* (2011)

Recombinaison homologue



Variation génétique et conséquence sur la protéine



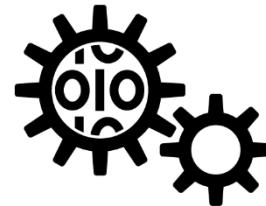
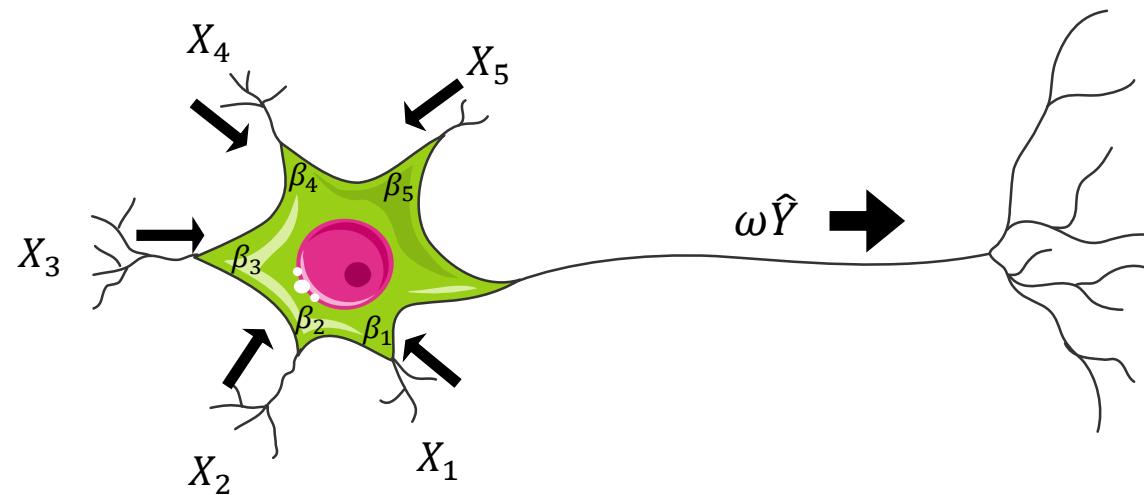
Intelligence artificielle



$$\hat{Y} = \beta_1 \times X_1 + \beta_2 \times X_2 + \cdots + \beta_p \times X_p$$

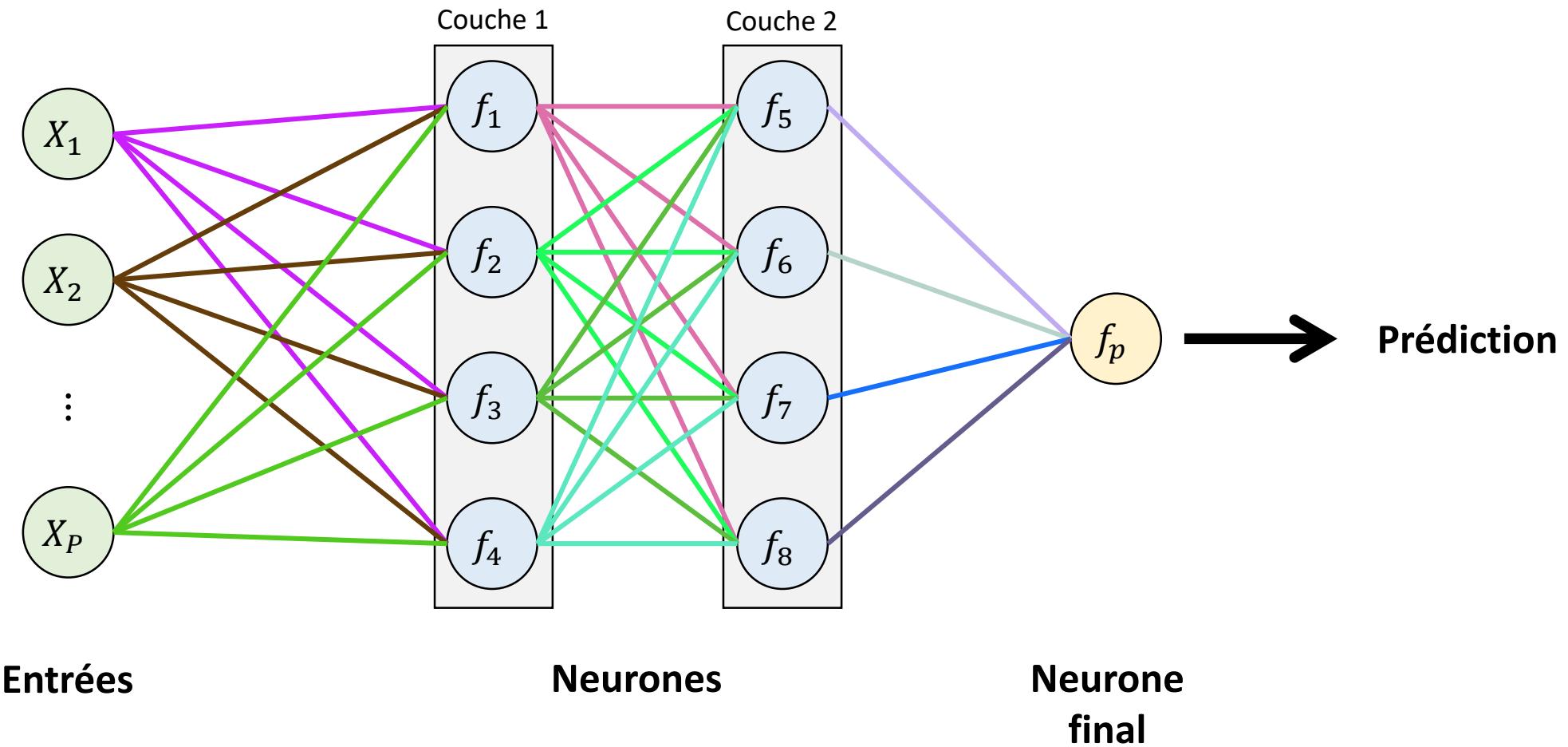
Machine learning

Deep learning



Machine learning

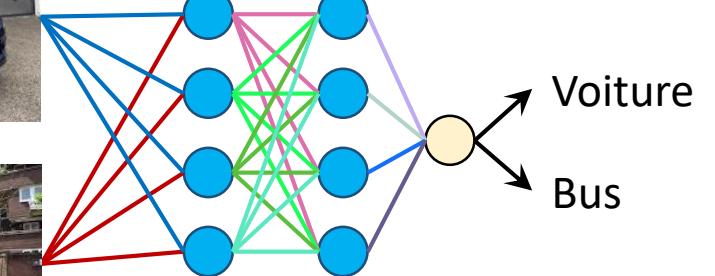
Les réseaux de neurones



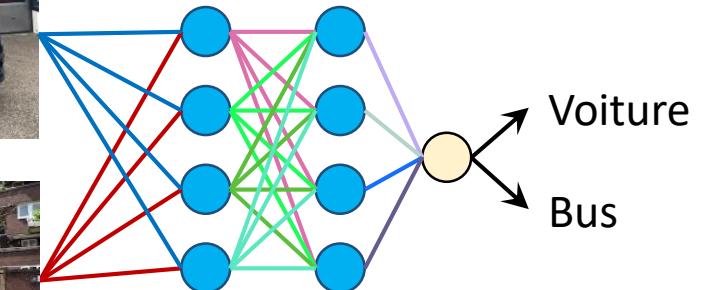
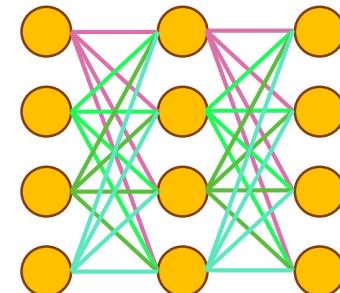
L'apport du *deep learning*



Machine learning

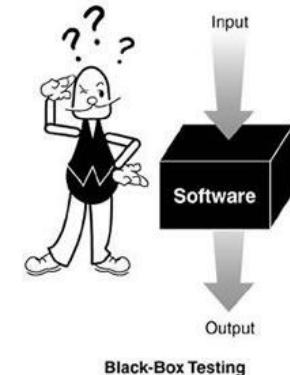
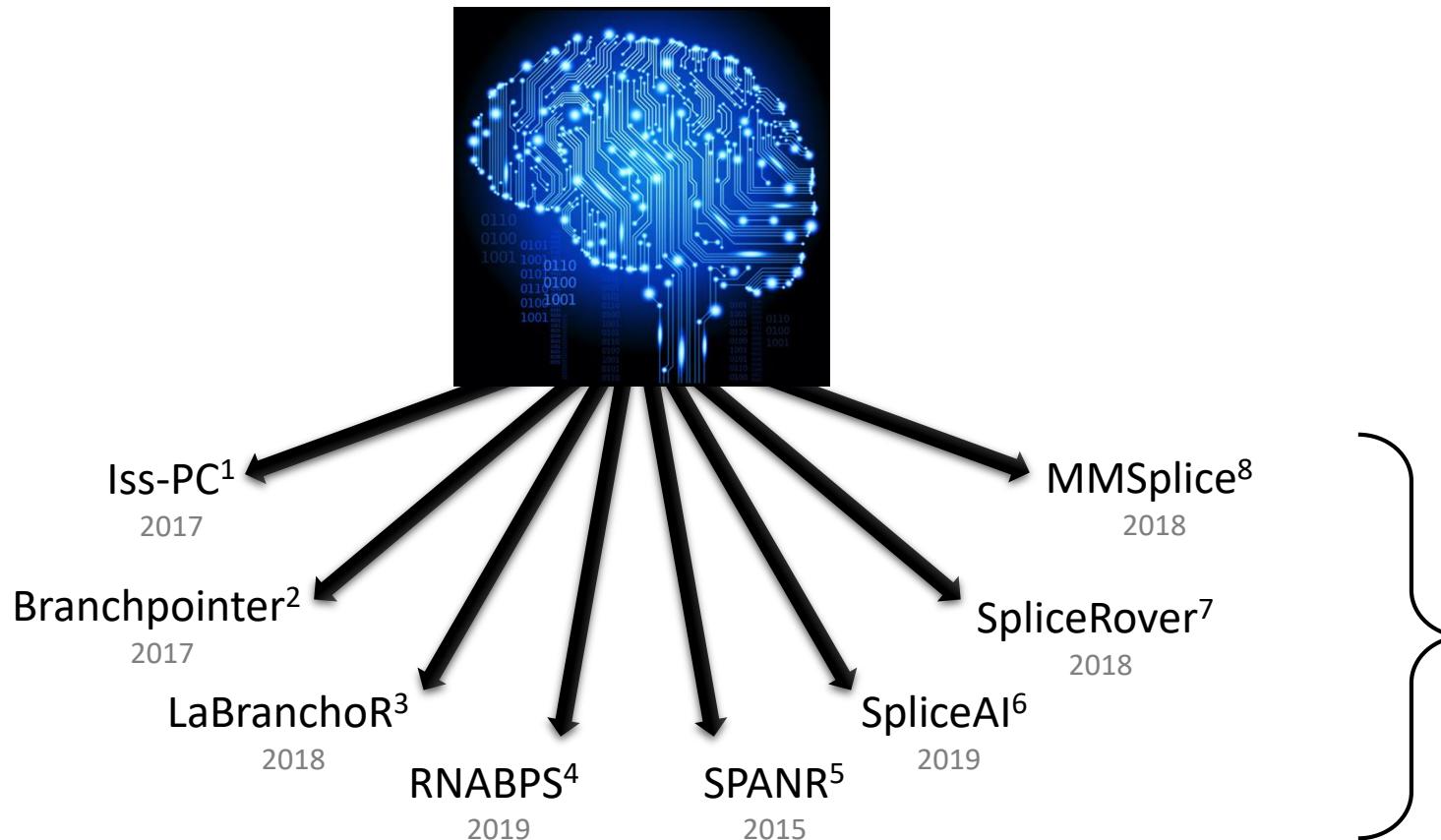


Deep learning



Deep learning et prédition épissage

Méthode en pleine essor



¹Xu *et al.*, *Sci. Rep.*, 2017

⁵Xiong *et al.*, *Science*, 2015

²Signal *et al.*, *Bioinformatics*, 2017

⁶Jaganathan *et al.*, *Cell*, 2019

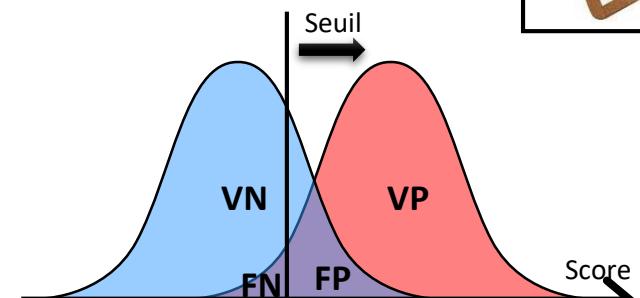
³Paggi *et al.*, *RNA*, 2018

⁷Zuallaert *et al.*, *Bioinformatics*, 2018

⁴Nazari *et al.*, *IEEE Access*, 2019

⁸Cheng *et al.*, *Genome Biol*, 2019

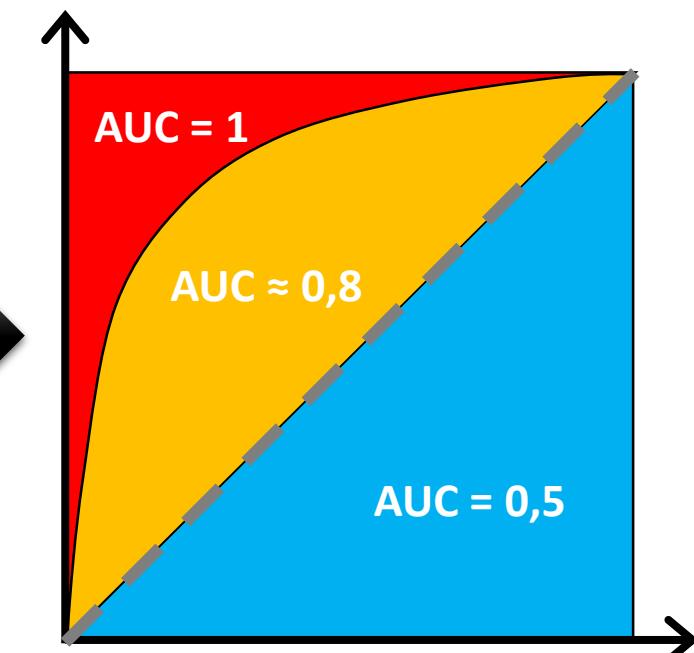
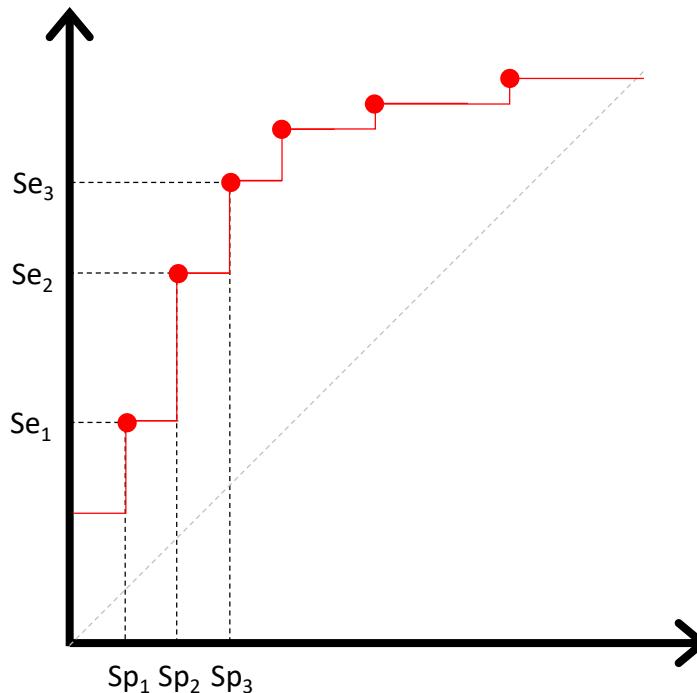
Évaluation des outils de prédition



Seuil 1		Observation	
Se_1/Sp_1		OUI	NON
Prédiction	OUI	VP	FP
		FN	VN

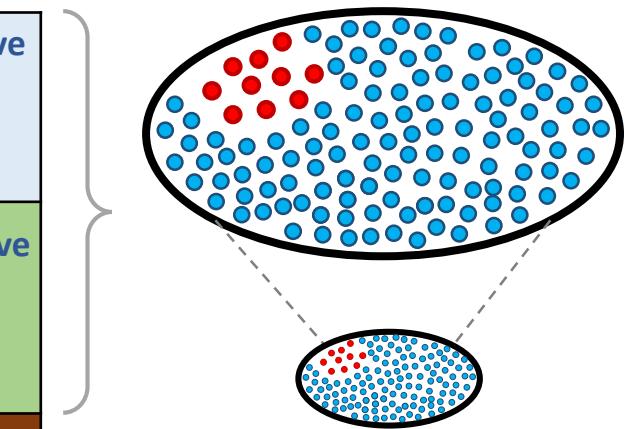
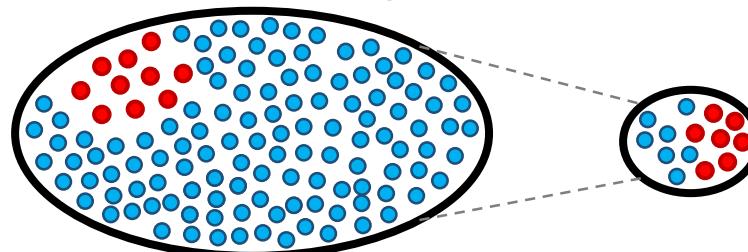
Seuil 2		Observation	
Se_2/Sp_2		OUI	NON
Prédiction	OUI	VP	FP
		FN	VN

Courbe ROC

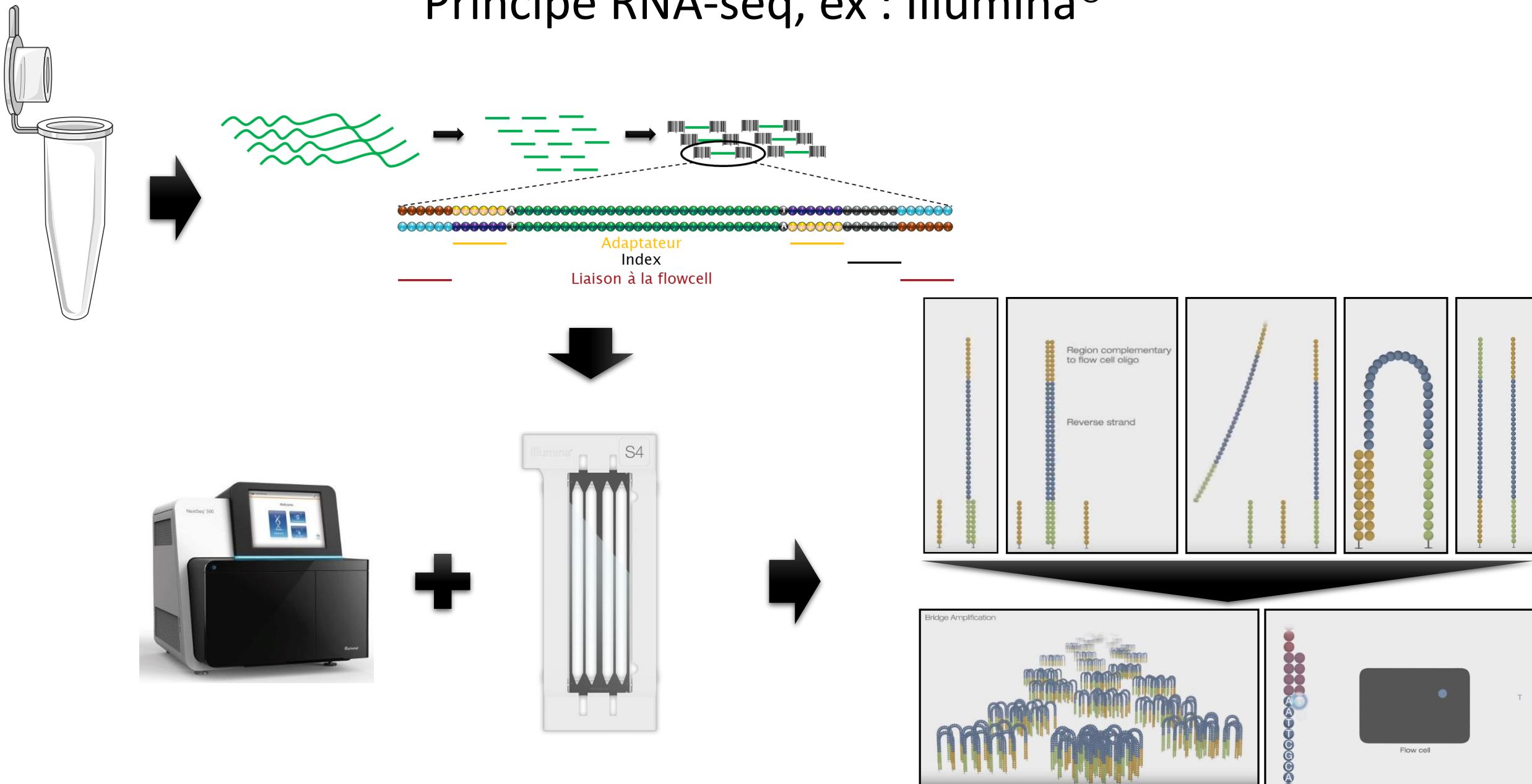


Évaluation des outils de prédiction

		Observation		
		Positive	Négative	
Prédition	Population totale			
	Positive	Vrai positif, VP	Faux positif, FP	Valeur prédictive positive (VPP), Précision $= \frac{VP}{VP + FP}$
	Négative	Faux négatif, FN	Vrai négatif, VN	Valeur prédictive négative (VPN) $= \frac{VN}{VN + FN}$
		<i>True positive rate (TPR), Recall, Sensibilité, Puissance</i> $= \frac{VP}{VP + FN}$	<i>Spécificité, Sélectivité, True negative rate (TNR)</i> $= \frac{VN}{VN + FP}$	Exactitude (accuracy) $= \frac{VP + VN}{VP + FP + VN + FN}$



Principe RNA-seq, ex : Illumina®

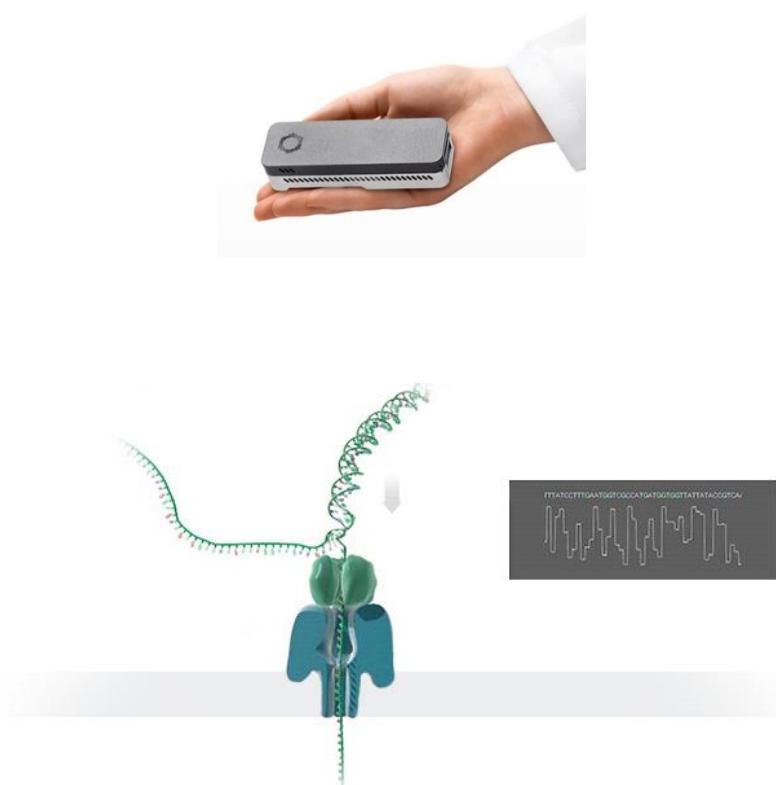


RNA-seq long read



Séquençage molécule unique

Transcrit de 1 à 100 kb



RNA-seq et analyses bioinformatiques

Alignment

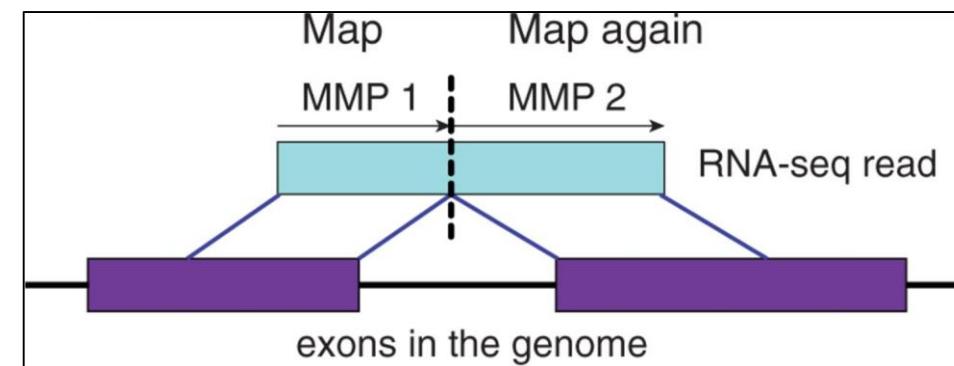
Annotation

Quantification

Ex : STAR*



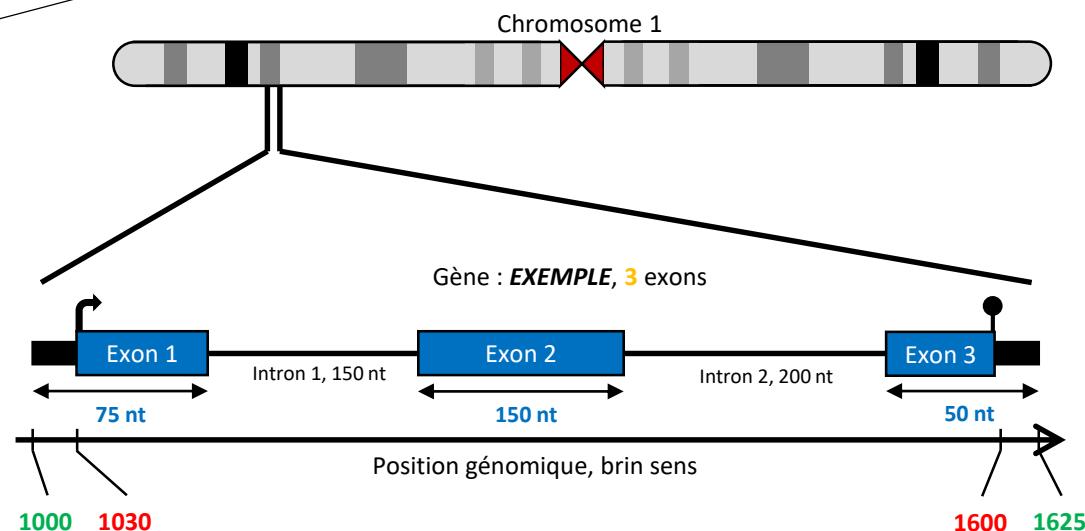
Génome



RNA-seq et analyses bioinformatiques



Ex : BEDtools*



Fichier BED (12 colonnes):

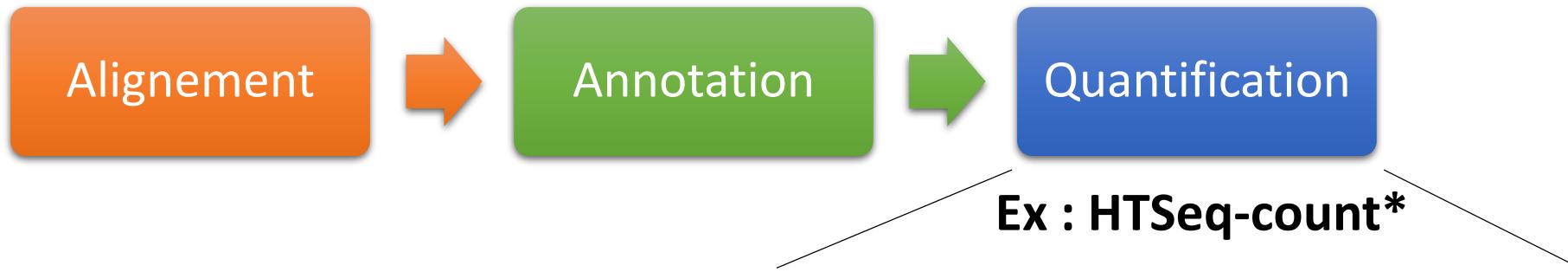
chr1	1000	1625	Exemple	0	+	1030	1600	0,0,255	3	75,150,50,	0,225,575,
------	------	------	---------	---	---	------	------	---------	---	------------	------------

Score entre 0 et 1000
Valeur définie par l'utilisateur

Code couleur RGB
0,0,255 = bleu

Position relative
Exon 1 : 0
Exon 2 : 75+150
Exon 3 : 75+150+150+200

RNA-seq et analyses bioinformatiques



	Ech_1	Ech_2	\cdots	Ech_j
Iso_1	C_{11}	C_{12}	\cdots	C_{1j}
Iso_2	C_{21}	C_{22}	\cdots	C_{2j}
\vdots	\vdots	\vdots	\ddots	\vdots
Iso_i	C_{i1}	C_{i2}	\cdots	C_{ij}

TO BE CONTINUED...

SpliceLauncher

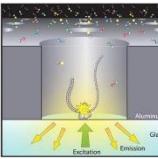


QC RNA-seq: towards guidelines for RNA-seq analyses (2/2)

BRCA1

IVS10-2A>G

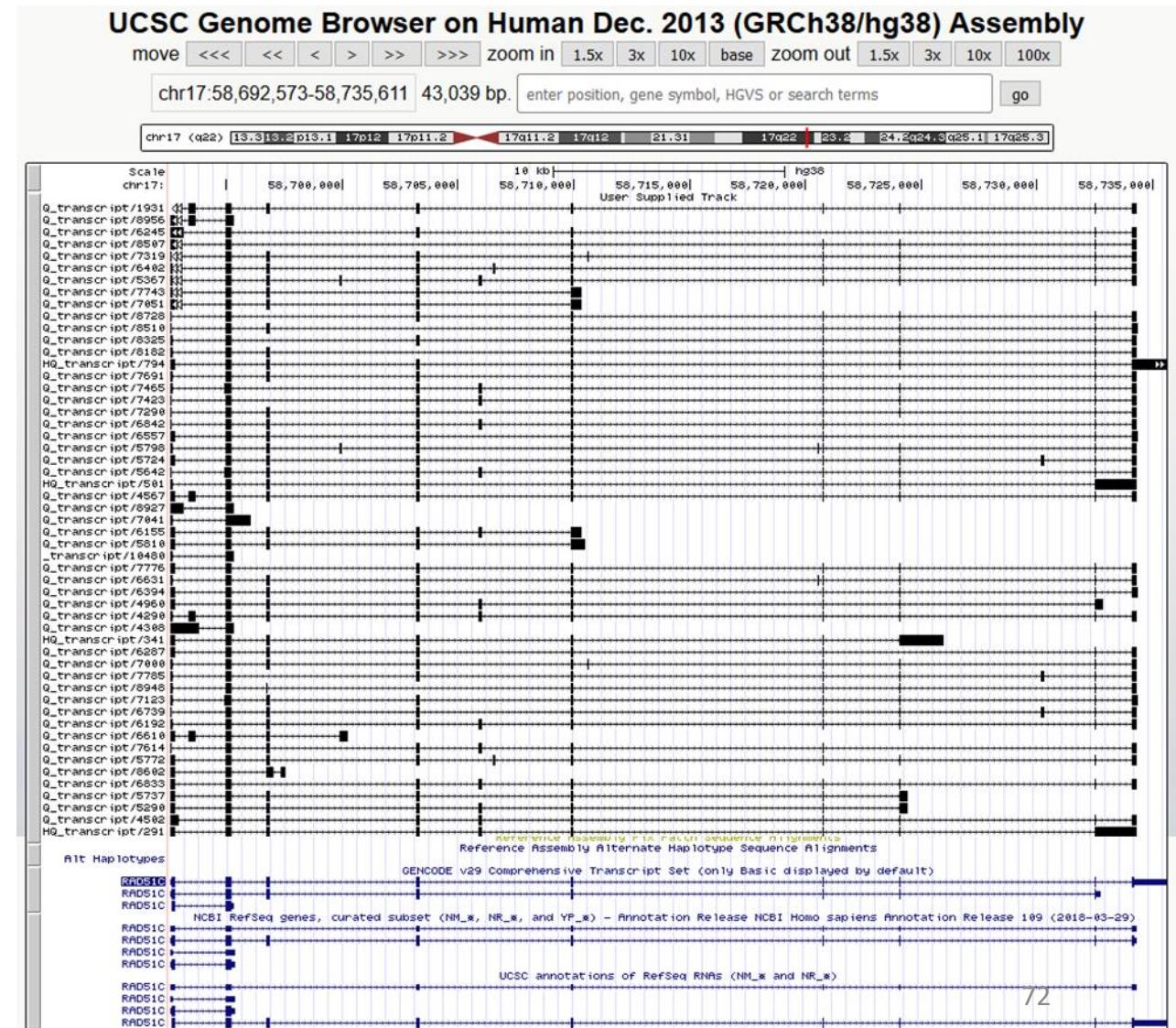




Long read sequencing for targeted RNA-seq (3/3)

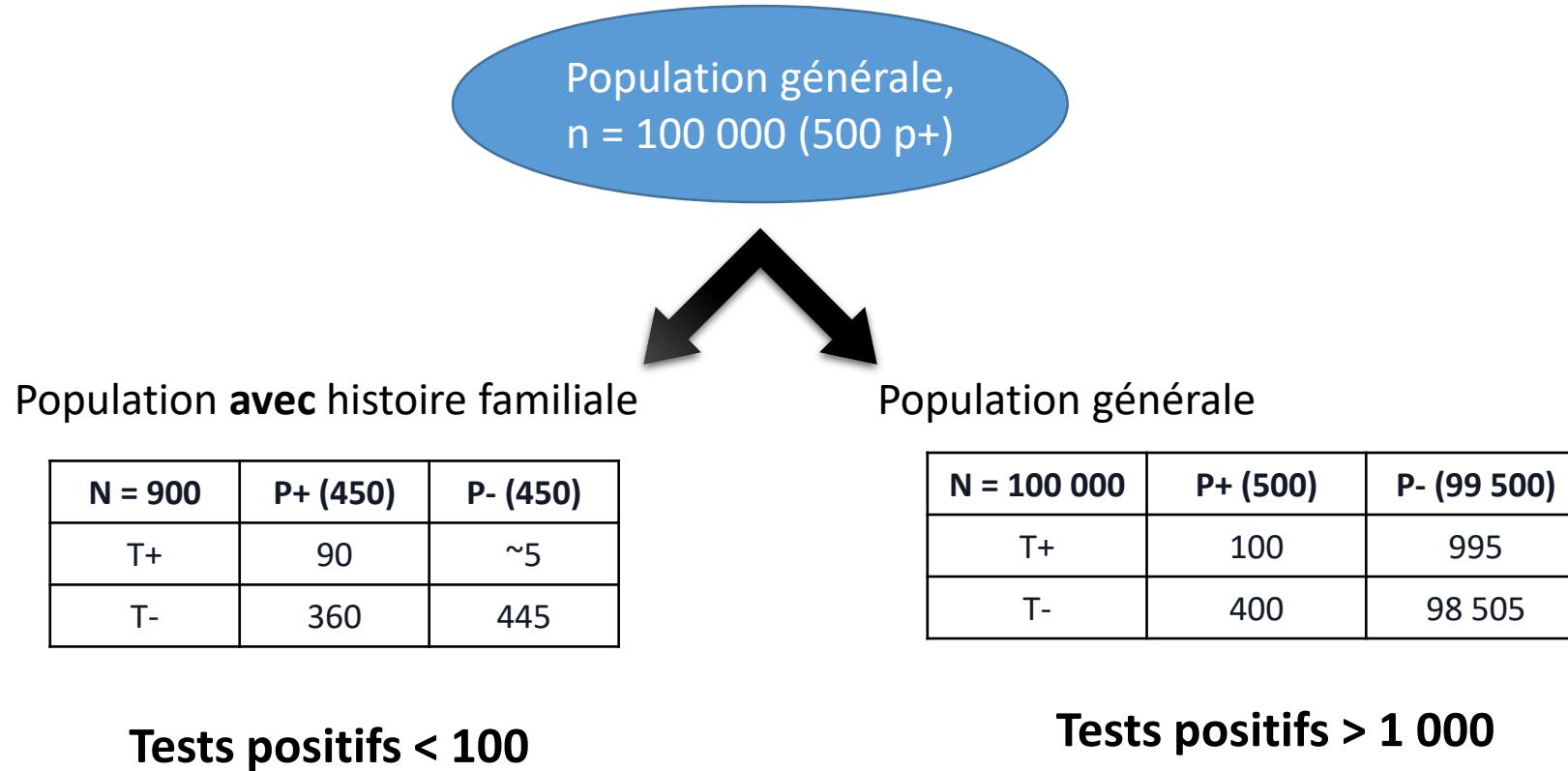
First results

Read count	<i>ACTB</i>	<i>BRCA1</i>	<i>BRCA2</i>	<i>PALB2</i>	<i>RAD51C</i>	<i>RAD51D</i>
<u>Without capture</u>	2 986	5	4	3	3	0
<u>With capture</u>	26	277	203	139	1107	14



Le piège des faux positifs

ACMG rules :
Se : 20 %
Sp : 99 %



Fausse conclusion

Faire un test génétique sur une population générale permet de trouver dix fois plus de prédisposition qu'un test selon les recommandations actuelles